

Requirements for Enterprise Search



December 2, 2004
© Stephen E. Arnold, 2004



Agenda for This Section

- Background
- Reputation analysis
- Wrap up

Questions at any time...

Why Requirements Get Ignored

- **Time pressure**
- Management demand
- Underestimate the challenges of search
- Misunderstand scaling / performance issues
- Assume that software can index, classify, and metatag

Rush Carefully



Why Requirements Get Ignored

- Time pressure
- **Management demand**
- Underestimate the challenges of search
- Misunderstand scaling / performance issues
- Assume that software can index, classify, and metatag

Dilbert Principle



Why Requirements Get Ignored

- Time pressure
- Management demand
- **Underestimate the challenges of search**
- Misunderstand scaling / performance issues
- Assume that software can index, classify, and metatag

Misjudge Task



Why Requirements Get Ignored

- Time pressure
- Management demand
- Underestimate the challenges of search
- Misunderstand scaling / performance issues
- **Assume that software can index, classify, and metatag**

Software... Not Always “Smart”



Why Requirements Are Important

- Search is complicated even with a “search toaster” from Google or Thunderstone
- Indexing and updates can slow a network so no other work can be done
- Human involvement is needed. Humans cost money
- Scaling is
 - Expensive
 - Slower than most people believe or accept

Google “Search Toaster”



GB-1001

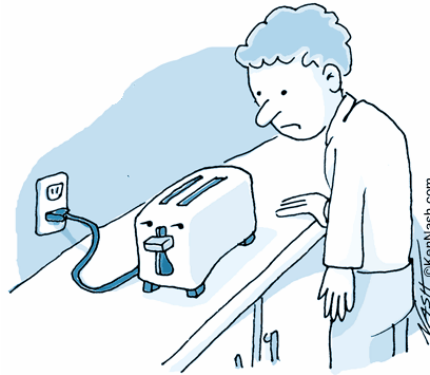


GB-5005



GB-8008

No Administration... Just Plug It In



The system administrator's friend... Plug in and go home.

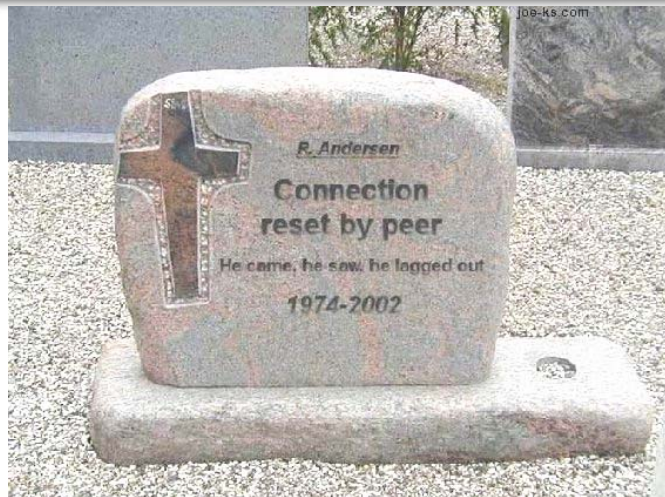
False Drops... No Hits...



Why Requirements Are Important

- Search is complicated even with a “search toaster” from Google or Thunderstone
- Indexing and updates can slow a network so no other work can be done
- Human involvement is needed. Humans cost money
- Scaling is
 - Expensive
 - Slower than most people believe or accept

Unresponsive Network



Why Requirements Are Important

- Search is complicated even with a “search toaster” from Google or Thunderstone
- Indexing and updates can slow a network so no other work can be done
- **Human involvement is needed. Humans cost money**
- Scaling is
 - Expensive
 - Slower than most people believe or accept

Humans Needed



Why Requirements Are Important

- Search is complicated even with a “search toaster” from Google or Thunderstone
- Indexing and updates can slow a network so no other work can be done
- Human involvement is needed. Humans cost money
- **Scaling is**
 - Expensive
 - Slower than most people believe or accept

Plan for a Big System



What Must Be Investigated?

- **Content that will be indexed?**
- Access to that content?
 - Security
 - Work flow
- What user interface?
 - Search box
 - Work flow (saved searches)
 - Yahoo-style facets
- Index update frequency?

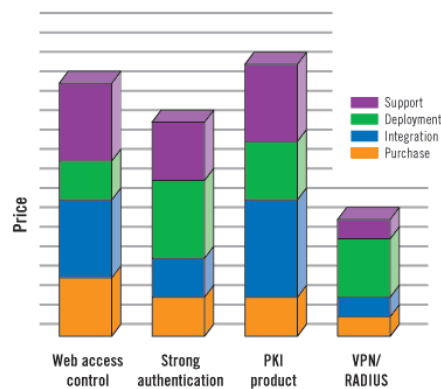
Humans Needed



What Must Be Investigated?

- Content that will be indexed?
- Access to that content?
 - Security
 - Work flow
- What user interface?
 - Search box
 - Work flow (saved searches)
 - Yahoo-style facets
- Index update frequency?

Log On

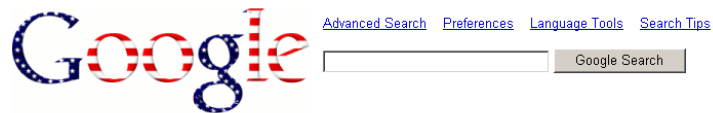


<http://www.securecomputing.com/index.cfm?skey=939>

What Must Be Investigated?

- Content that will be indexed?
- Access to that content?
 - Security
 - Work flow
- What user interface?
 - Search box
 - Work flow (saved searches)
 - Yahoo-style facets
- Index update frequency?

Google



Yahoo



What Must Be Investigated?

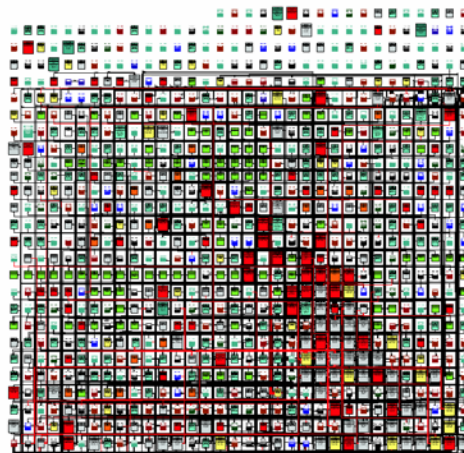
- Content that will be indexed?
- Access to that content?
 - Security
 - Work flow
- What user interface?
 - Search box
 - Work flow (saved searches)
 - Yahoo-style facets
- Index update frequency?

Update One Hour after Upset

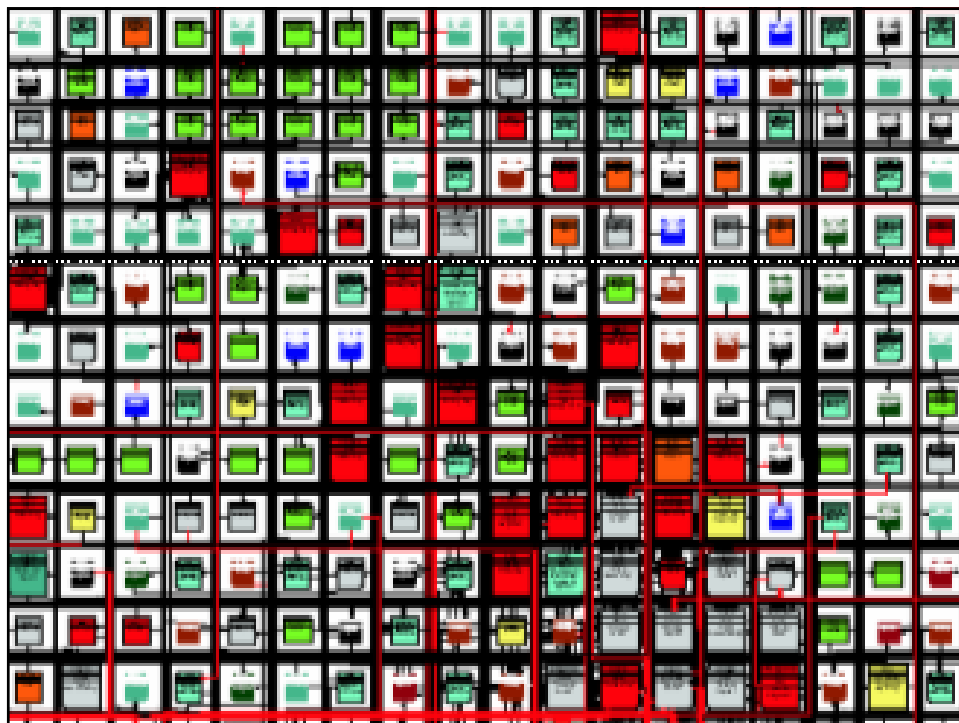
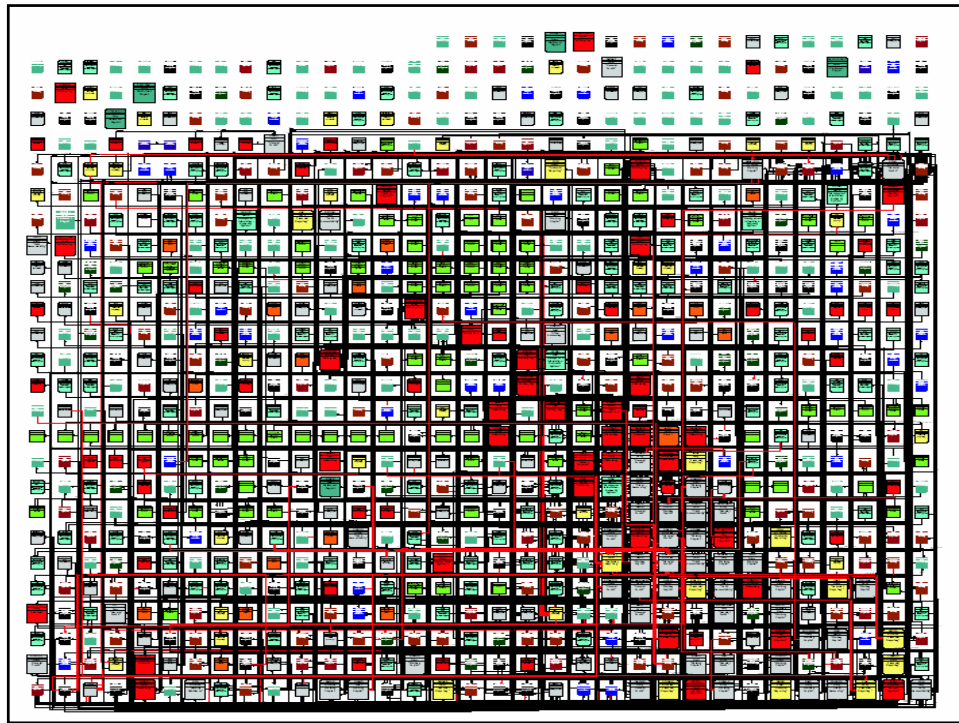


No update within one hour of Patriots' loss.

Red=Sensitive, Green=Not Sensitive, Other=?



Source: Catherine Santana, Director FMMP, DoD, 2002



What Must Be Investigated?

- **Who can access the system?**
- What usage tracking is needed?
- What are the file types to be indexed?
- What legacy systems' content must be indexed?
- Will source documents be served from:
 - A data warehouse
 - The server or machine where documents reside?

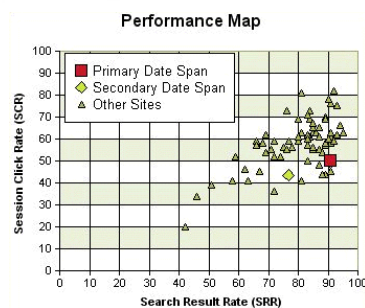
Jobs, Not Individuals



What Must Be Investigated?

- Who can access the system?
- **What usage tracking is needed?**
- What are the file types to be indexed?
- What legacy systems' content must be indexed?
- Will source documents be served from:
 - A data warehouse
 - The server or machine where documents reside?

Metrics



Most-Recent Sessions

This page lists the most recent user sessions. Each session is identified either by the IP number or (when available) the domain name of the searcher.

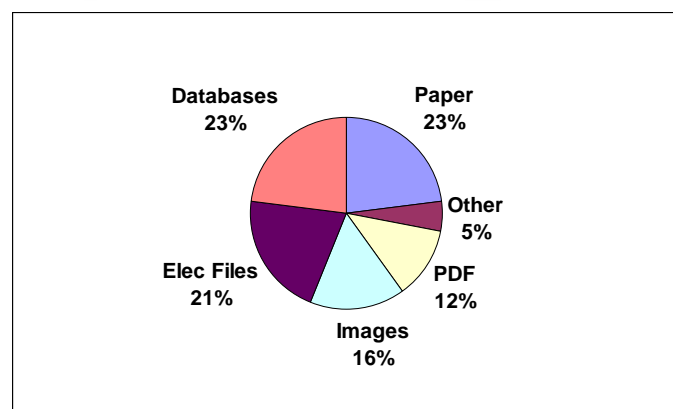
Number of rows: 10 [Apply](#) [Print Table](#) [Export Table](#)

Primary Date Span April 30 2003 - May 13 2003		Secondary Date Span January 30 2003 - April 29 2003	
#	Session	#	Session
01	...adsl.prowd.net	5/13/2003	...thruvonder.co.uk
02	...bluepounder.co.uk	5/13/2003	...72.cable.rcf.com
03	...95.6sl.papers.com	5/13/2003	...615-2.Nelond.nl
04	...71.6sl.papers.com	5/13/2003	...78.6sl.papers.com
05	...interbusiness.it	5/13/2003	...79.6sl.papers.com
06	...61.6.208.10	5/13/2003	...167.mn.88a.com
07	...8.abo.wanadoo.fr	5/13/2003	...82.6sl.psl.co.uk
08	...8.abo.wanadoo.fr	5/13/2003	...cl.dsl.psl.co.uk
09	...8.abo.wanadoo.fr	5/13/2003	...cl.dsl.psl.co.uk
10	...8.abo.wanadoo.fr	5/13/2003	...interbusiness.it

What Must Be Investigated?

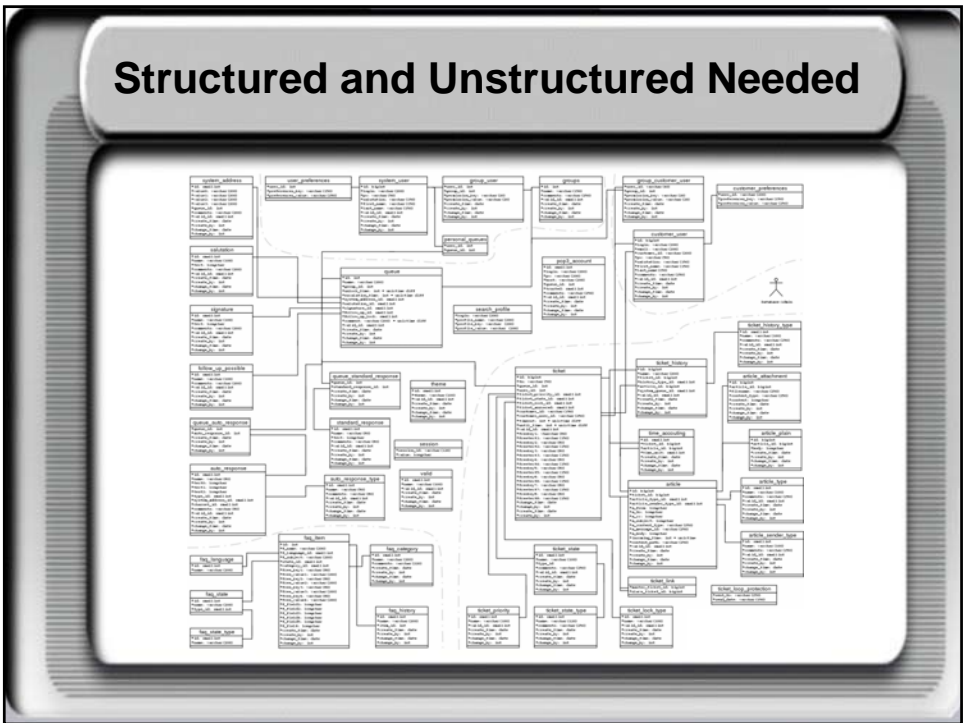
- Who can access the system?
- What usage tracking is needed?
- **What are the file types to be indexed?**
- What legacy systems' content must be indexed?
- Will source documents be served from:
 - A data warehouse
 - The server or machine where documents reside?

Document Formats



Source: AMR Research, September 2004

Databased Information... High Value



What Must Be Investigated?

- Who can access the system?
- What usage tracking is needed?
- What are the file types to be indexed?
- **What legacy systems' content must be indexed?**
- Will source documents be served from:
 - A data warehouse
 - The server or machine where documents reside?

Legacy Systems

```
x3270-2 ibmlink.advantis.com
```

FileOptions

SVM0201P
SYSTEM: NAM0SM94
TERMID: IBM0TSP
CUSTOMER ASSISTANCE: PLEASE CALL 1-800-727-2222
DATE: 08/07/17
TIME: 17:21:23

Welcome to

===== (R)
===== **
===== ** *
===== ** *
===== ** *
===== ** *
===== ** *
===== ** *
===== ** *
===== ** *
===== ** *
===== ** *
===== ** *
===== ** *
===== ** *

(R) Registered trademark of the IBM Corporation
(C) Copyright International Business Machines Corporation 1985, 1993

ACCOUNT... USERID... PASSWORD...
Enter desired product or service, or press the HELP key (PF1) for assistance.
====>

021/01

What Must Be Investigated?

- Who can access the system?
- What usage tracking is needed?
- What are the file types to be indexed?
- What legacy systems' content must be indexed?
- Will source documents be served from:
 - A data warehouse
 - The server or machine where documents reside?

Storage, Not a Single Hard Drive



Challenges

- Data waves, not floods
- Requirements versus budget
- “Nice to have” versus “must have”
- Knowledge versus assumptions
- Experience versus inexperience
- Time versus expectations
- Vendor’s promises versus system realities
- Existing infrastructure versus search systems’ storage and computational hunger

Data: Waves, not Floods



- Large volumes of data
- Fast moving
- Changes
 - Quite important
 - Major challenge

Requirements Pitfalls

- **Users explain search in terms of Google, not their specific work-related needs**
- Customization
- Resources limited; therefore, risk increases
- Too little research; therefore, decisions are guesses
- Ignore / underestimate:
 - Copyright issues
 - Access procedures
 - Indexing issues

The Google Effect



Requirements Pitfalls

- Users explain search in terms of Google, not their specific work-related needs
- **Customization**
- Resources limited; therefore, risk increases
- Too little research; therefore, decisions are guesses
- Ignore / underestimate:
 - Copyright issues
 - Access procedures
 - Indexing issues

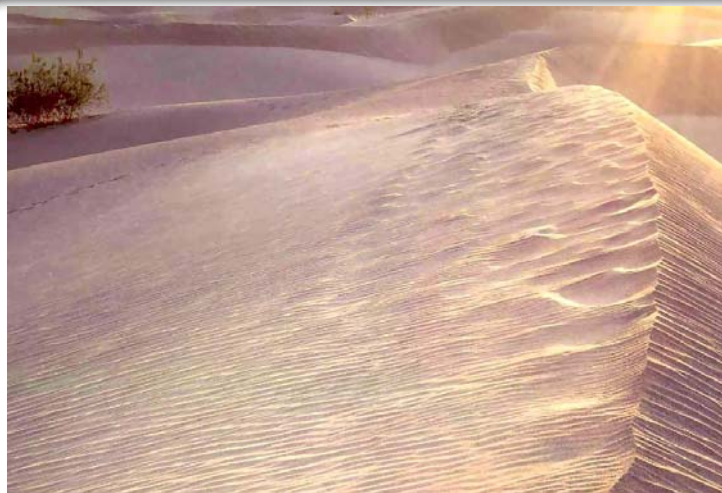
Customization Costs



Requirements Pitfalls

- Users explain search in terms of Google, not their specific work-related needs
- Customization
- **Resources limited; therefore, risk increases**
- Too little research; therefore, decisions are guesses
- Ignore / underestimate:
 - Copyright issues
 - Access procedures
 - Indexing issues

Adequate Resources... Essential



Requirements Pitfalls

- Users explain search in terms of Google, not their specific work-related needs
- Customization
- Resources limited; therefore, risk increases
- **Too little research; therefore, decisions are guesses**
- Ignore / underestimate:
 - Copyright issues
 - Access procedures
 - Indexing issues

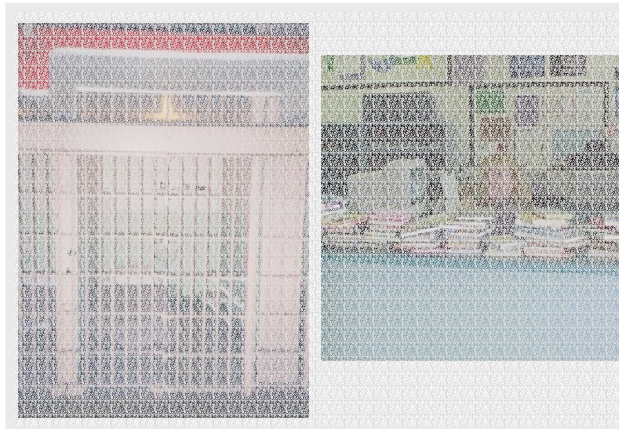
Guesses Not Good



Requirements Pitfalls

- Users explain search in terms of Google, not their specific work-related needs
- Customization
- Resources limited; therefore, risk increases
- Too little research; therefore, decisions are guesses
- Ignore / underestimate:
 - Copyright issues
 - Access procedures
 - Indexing issues

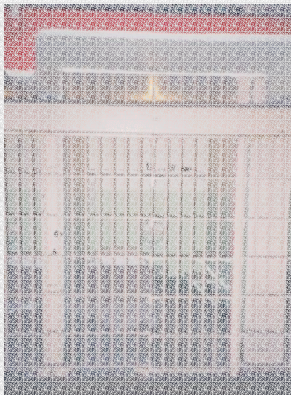
Copyright... Access... Manual Work



Copyright... Access... Manual Work



Copyright... Access... Manual Work



Lessons Learned

- **No free.lunch... search can and will top \$1 million in its first year in mid-sized organizations**
- Fruitcake data... data harmonization requires manual effort plus software
- Silos of data... update and security issues persist
- Costs go up faster than inflation; for example, 212 percent over an 18 month period
- Technical problems can be expensive to solve

The Key to Rocket Science



Lessons Learned

- No freelunch... search will easily top \$1 million in its first year in mid-sized organizations
- **Fruitcake data... data harmonization requires manual effort plus software**
- Silos of data... update and security issues persist
- Costs go up faster than inflation; for example, 212 percent over an 18 month period
- Technical problems can be expensive to solve

Data Fruitcake



Lessons Learned

- No freelunch... search can easily top \$1 million in its first year in mid-sized organizations
- Fruitcake data... data harmonization requires manual effort plus software
- **Silos of data... update and security issues persist**
- Costs go up faster than inflation; for example, 212 percent over an 18 month period
- Technical problems can be expensive to solve

One Data Silo Means...



More Silos... Not Integration



Lessons Learned

- No freelunch... search can easily top \$1 million in its first year in mid-sized organizations
- Fruitcake data... data harmonization requires manual effort plus software
- Silos of data... update and security issues persist
- **Costs go up faster than inflation; for example, 212 percent over an 18 month period**
- Technical problems can be expensive to solve

Costs Can Explode



Nomar / A-Rod / Jeter

Lessons Learned

- No freelunch... search can easily top \$1 million in its first year in mid-sized organizations
- Fruitcake data... data harmonization requires manual effort plus software
- Silos of data... update and security issues persist
- Costs go up faster than inflation; for example, 212 percent over an 18 month period
- Technical problems can be expensive to solve

Complicated Problems

$$(1): \frac{\partial C_i}{\partial t} + \frac{\partial J_x}{\partial x} + \frac{\partial J_y}{\partial y} + \frac{\partial J_z}{\partial z} = q_{, i=0, w, g}$$

The fluxes:

$$(\vec{J})_0 = \frac{\rho_{osc}}{B_0} \vec{V}_0$$

$$(2): (\vec{J})_w = \frac{\rho_{WSC}}{B_w} \vec{V}_w$$

$$(\vec{J})_g = \frac{\rho_{gsc}}{B_g} \vec{V}_g + \frac{R_{so} \rho_{gsc}}{B_0} \vec{V}_0 + \frac{R_{sw} \rho_{gsc}}{B_w} \vec{V}_w$$

Requirements Tips

- **Form a search “team”. Work as a team.**
- Identify the key stakeholders. Discuss their search needs.
- Use a Web survey for general information.
- Communicate to keep expectations realistic.
- If search will be embedded in work flow, concentrate data collection where search touches existing processes.
- If resources are available, consider outside requirements specialists.

A Team ... Search Has Many Parts



Requirements Tips

- Form a search “team”. Work as a team.
- **Identify the key stakeholders. Discuss their search needs.**
- Use a Web survey for general information.
- Communicate to keep expectations realistic.
- If search will be embedded in work flow, concentrate data collection where search touches existing processes.
- If resources are available, consider outside requirements specialists.

Find the Stakeholders



Requirements Tips

- Form a search “team”. Work as a team.
- Identify the key stakeholders. Discuss their search needs.
- **Use a Web survey for general information.**
- Communicate to keep expectations realistic.
- If search will be embedded in work flow, concentrate data collection where search touches existing processes.
- If resources are available, consider outside requirements specialists.

Use a Web Survey

Do you access the service from your home?

☒ No
☐ Yes

For member offices only. Does anyone in the home state office use the Senate News Wire?

☒ No
☐ Yes. If yes, how many and what are their titles?

How do you like to search for information? (check all that apply)

☐ Personal Views
☐ Point and click on subject areas in Shared Views
☐ Manual search (search tab)
☐ Other (You must type an answer if you select this choice)

Do you at times limit your search to one or more specific sources?

☒ No
☐ Yes
☐ Did not know I could

Requirements Tips

- Form a search “team”. Work as a team.
- Identify the key stakeholders. Discuss their search needs.
- Use a Web survey for general information.
- **Communicate to keep expectations realistic.**
- If search will be embedded in work flow, concentrate data collection where search touches existing processes.
- If resources are available, consider outside requirements specialists.

Search... No Pot of Gold



Requirements Tips

- Form a search “team”. Work as a team.
- Identify the key stakeholders. Discuss their search needs.
- Use a Web survey for general information.
- Communicate to keep expectations realistic.
- **If search will be embedded in work flow, concentrate data collection where search touches existing processes.**
- If resources are available, consider outside requirements specialists.

Focus Where Search Hits Work Flow



Requirements Tips

- Form a search “team”. Work as a team.
- Identify the key stakeholders. Discuss their search needs.
- Use a Web survey for general information.
- Communicate to keep expectations realistic.
- If search will be embedded in work flow, concentrate data collection where search touches existing processes.
- **If resources are available, consider outside requirements specialists.**

Consider Consultants... Save Time



Tomorrow...the “Semantic Web”?

