



Search Engine Technologies

Status and Directon

Ingvar.Aaberg@fast.no
February 1, 2006



Outline

- **Background and introduction**
- **Anatomy of a search engine**
- **Some thoughts on the web search market**
- **Enabling technologies**

FAST Background

Corporate Overview

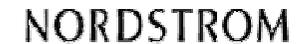
- **Leadership**

- Founded in 1997
- Public company (OSE: 'FAST')
- Profitable and well capitalized
- Revenue growth = 50%
- > 3,000 customers
- 40+ PhDs



- **Focus on Enterprise Search**

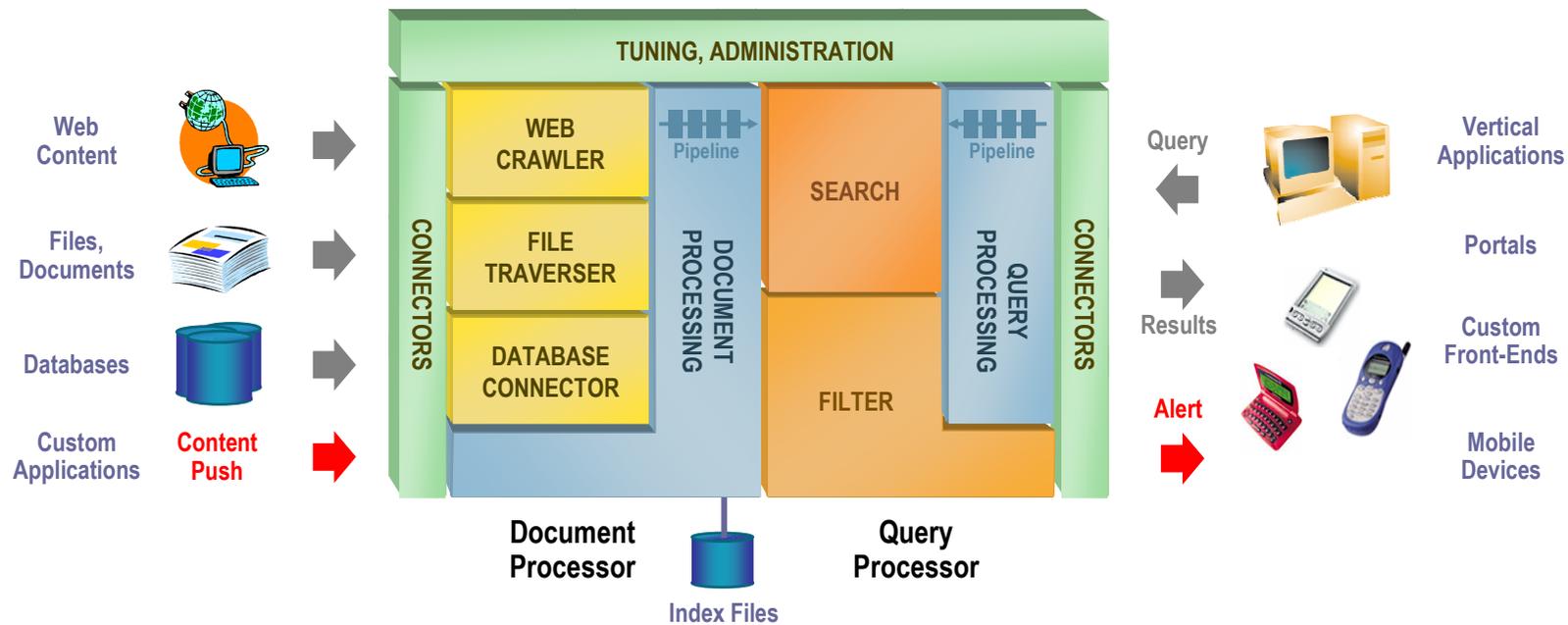
- Sold internet business (alltheweb.com)
- Acquiring enterprise search complements
- Ranked market leader by Gartner



Search Engine

How It Works

- **Two asynchronous independent components**
 - Document processor
 - Query processor



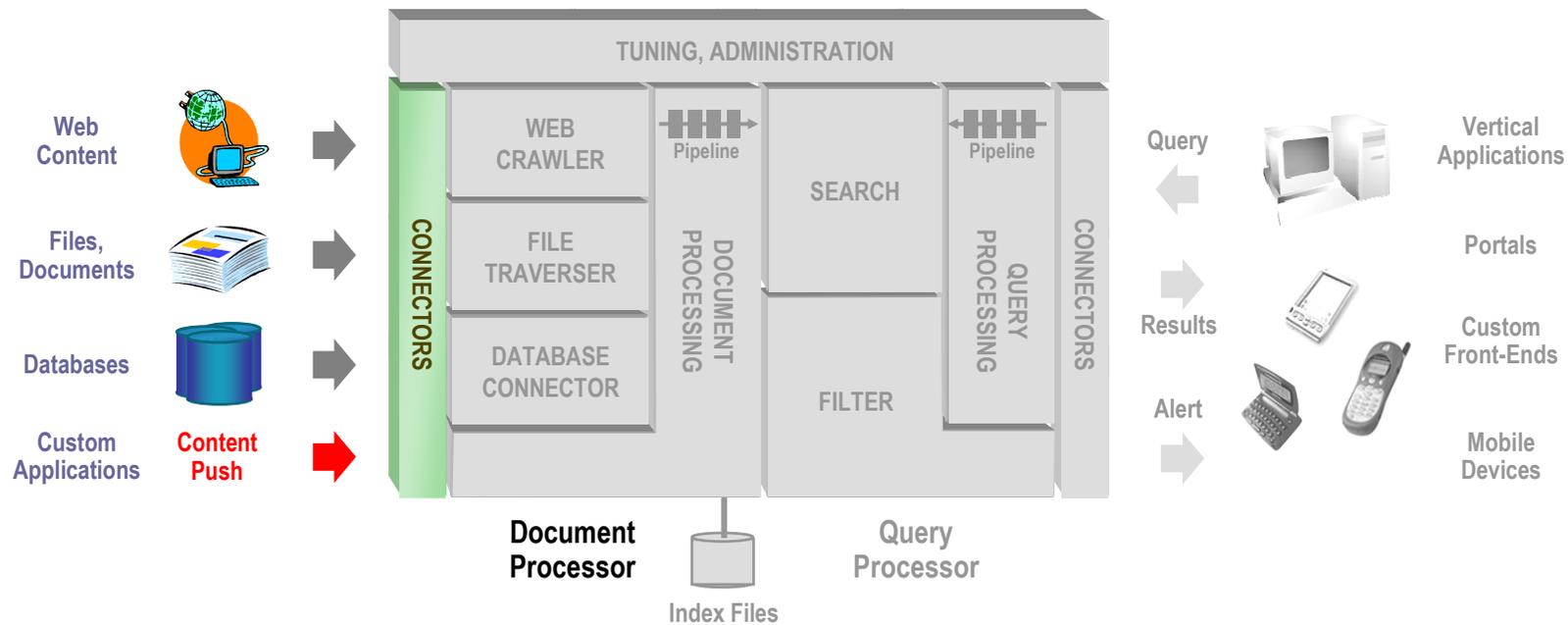
Search Engine

How It Works

1

Connect to content sources

- Structured (RDBMS)
- Unstructured (Web, Office, etc.)
- Semi-structured (XML)

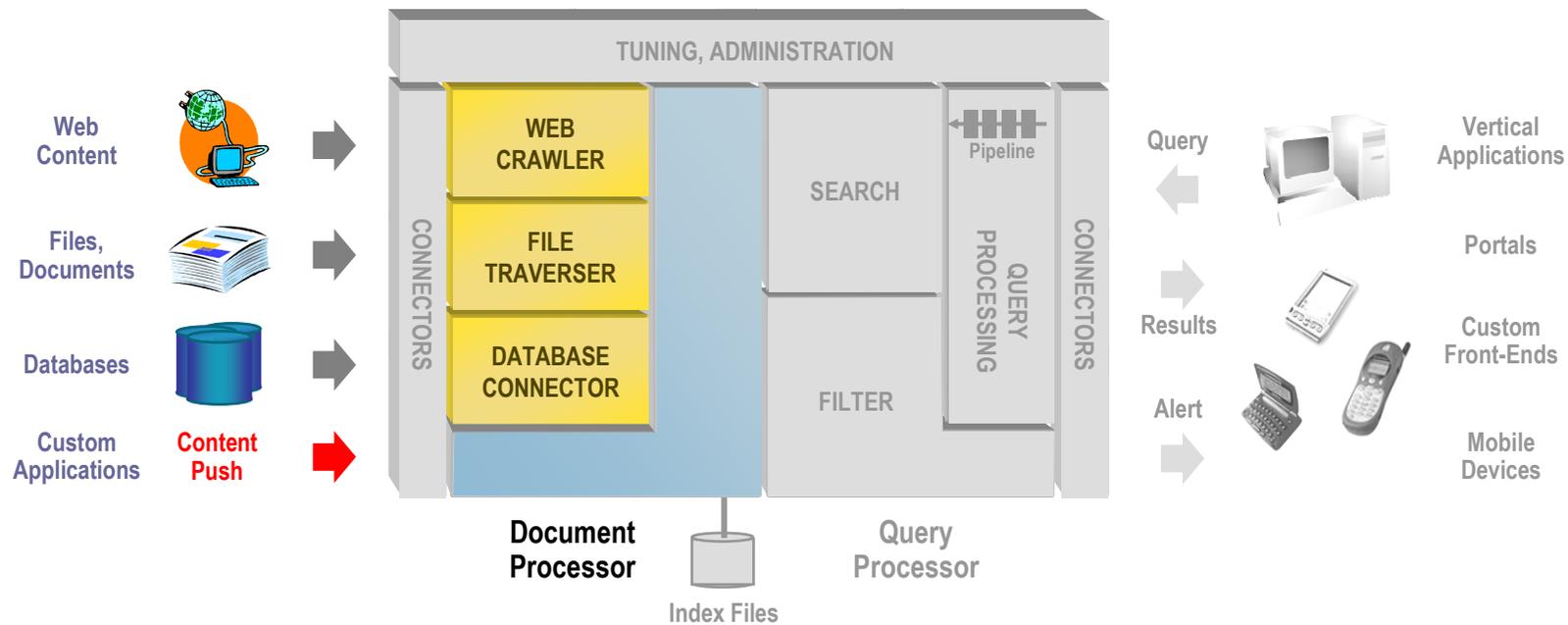


Search Engine

How It Works

2 Get content data

- Gather data wholesale by incrementally crawling links
- Push changed data content directly to processing stage



Document processing

- **Processes the content before it gets indexed**
 - Documents flow through a pipeline of processing stages
 - Highly customizable

- **Example processing stages**

- Format, language and encoding detection pdf; english; iso-8859-1
- Format and encoding normalization pdf → html; iso-8859-1 → utf-8
- HTML parsing
- Entity extraction venus williams; arthur ashe; ...
- Vectorization {(venus williams, 1.0), (wimbledon, 0.81), (center court, 0.65), ...}
- Categorization sports/tennis
- Lemmatization and synonym handling mouse ~ mice; car ~ automobile
- Anchor text harvesting
- ...

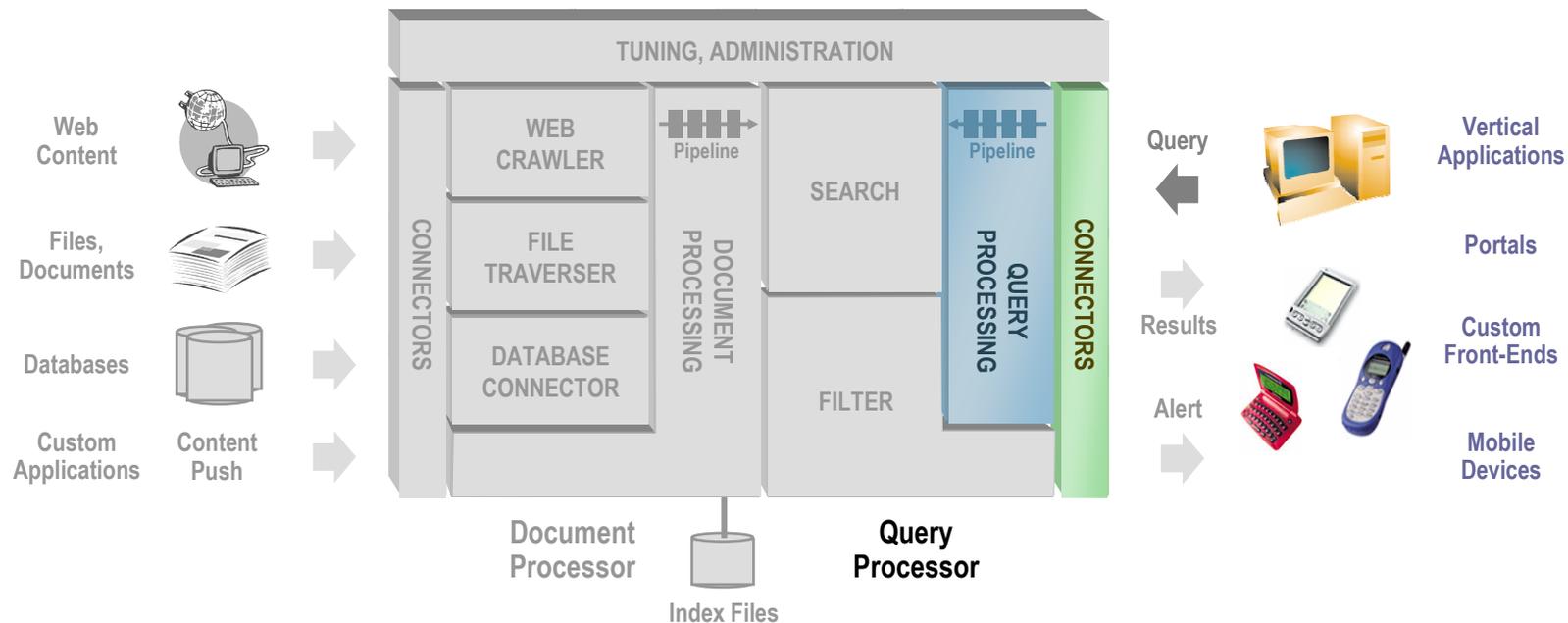
Search Engine

How It Works

4

Interpret query

- Use query language or query API
- Convert and process query through query pipeline:
 - Linguistic processing
 - Custom logic (e.g. query term modification/addition)



Content Enrichment

- **Processes the query before it gets sent to the search engine**

- Queries flow through a pipeline of processing stages
- Highly customizable

- **Example processing stages**

- Stopword handling
- Phrasing and antiphrasing
- Spellchecking
- Natural language handling
- Lemmatization and synonym handling
- Ontologies
- ...

red cross → "red cross"; where can i find information about cars → cars

brittany speers → britney spears

television under 200 dollars → television AND price:<200

mouse ~ mice; car ~ automobile

xsara < citröen < car < vehicle

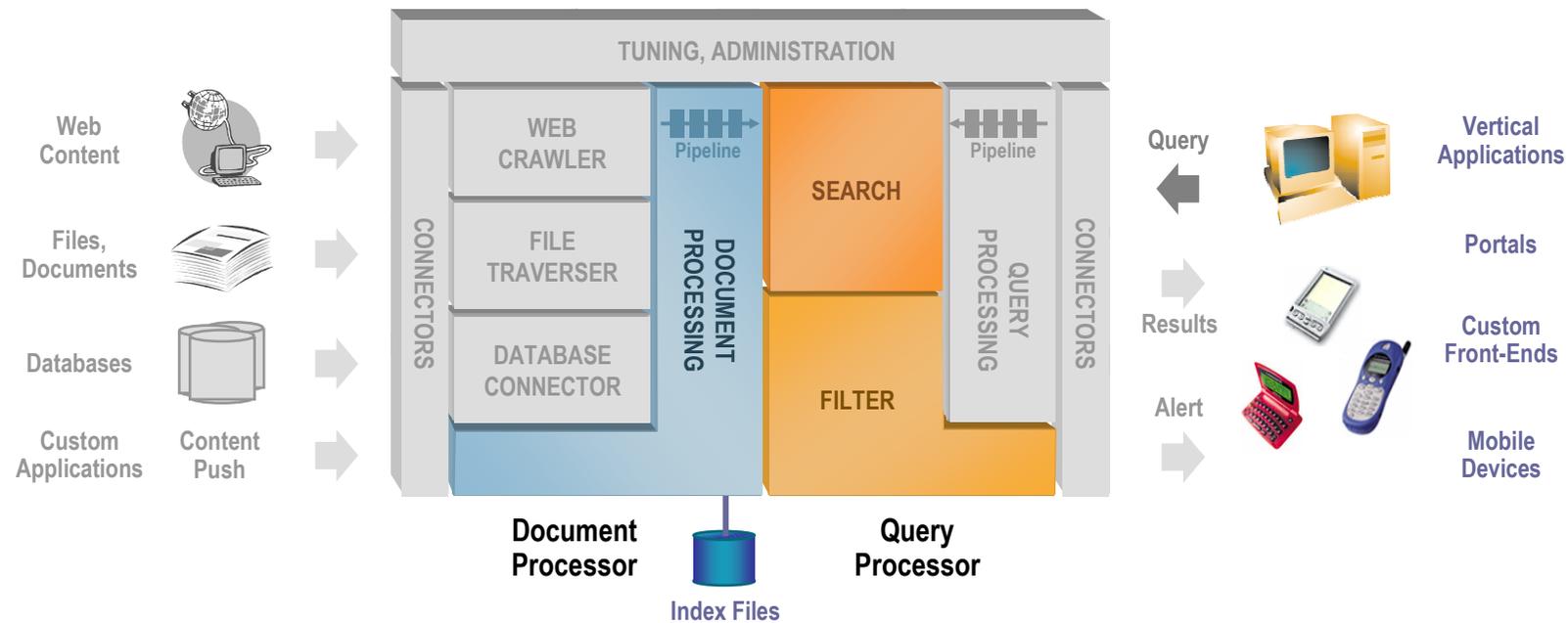
Search Engine

How It Works

5

Get results

- Get crawled results from index files
- Get live results directly from processor



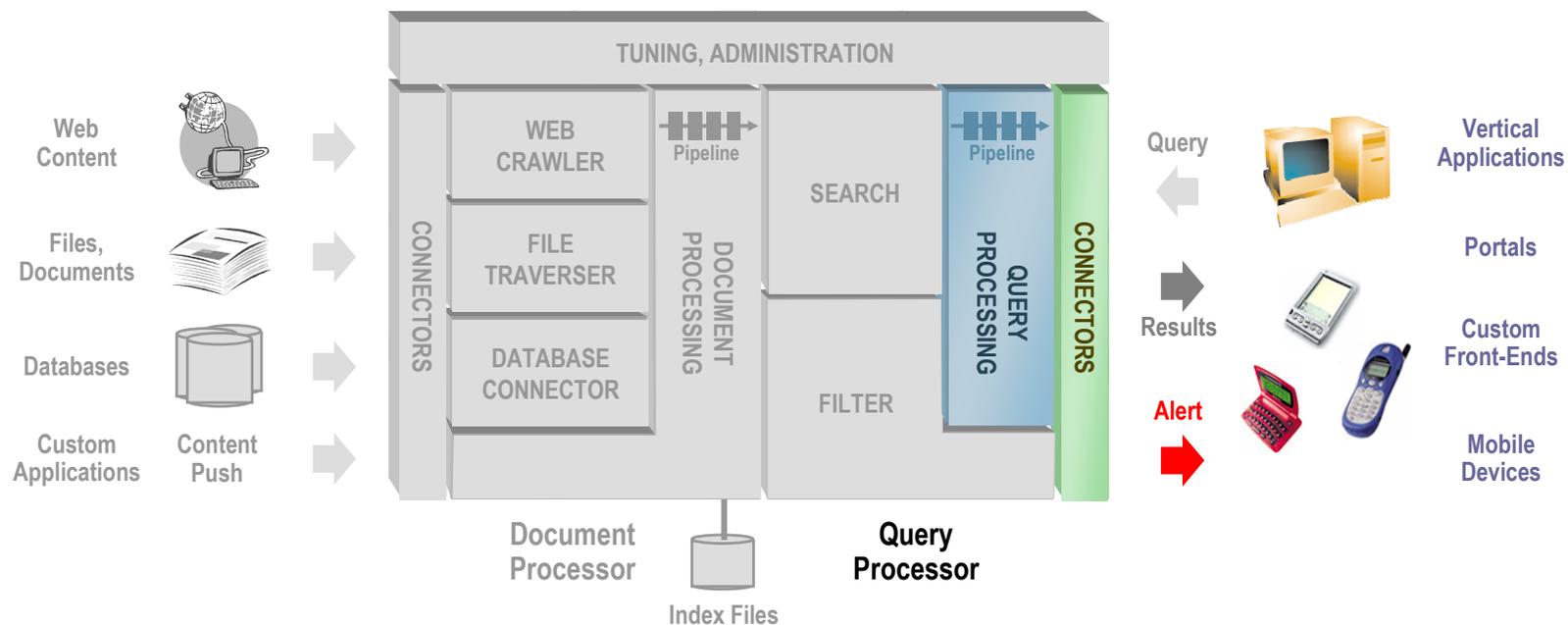
Search Engine

How It Works

6

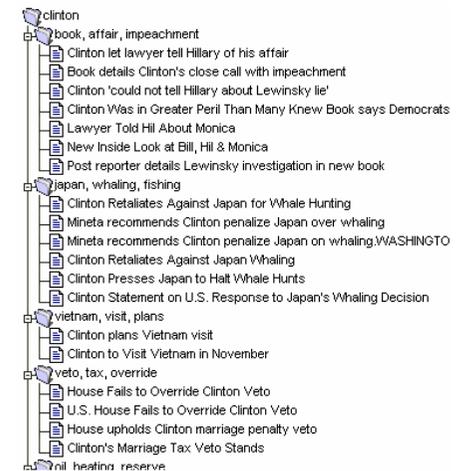
Return results to user

- Convert and process results through query pipeline:
 - Resort, filter for security, organize for dynamic drilldown
- Pass results on to application (generated or through API)
- Push results to alert engine and then external environment (e.g. mail, queue)



Result processing

- **Unsupervised clustering**
 - Analyzes on-the-fly similarities across matching documents
 - Group together documents having similar content
 - Provides a bird's-eye view of topical spread
- **Query refinement suggestions**
 - Examines the distribution of meta data
 - Builds a histogram of values
 - Provides a means for slicing and dicing the result set
- **Filtering**
 - E.g., dynamic duplicate detection



Category	Item	Count
Emails	jltke@ap.org	3
	avire@marion.gannett.com	2
	tree@ap.org	2
	jmongeon22@hotmail.com	1
	mark@heritagesc.com	1
	mcnall@together.net	1
	meshell@ic.netcom.com	1
	millik	1
	nynar	1
	parise	1
People	Roger Federer	68
	Andy Roddick	51
	Lindsay Davenport	50
	Andre Agassi	48
	Maria Sharapova	45
	Serena Williams	45
	Alicia Molik	36
	Marat Safin	34
	Nikolay Davydenko	25
	Joachim Johansson	25
Keywords	australian open	120
	melbourne	63
	roger federer	50
	open	48
	grand slam	44
	sydney	39
	wimbledon	37
	andy roddick	37
	hevit	36
	andre agassi	36

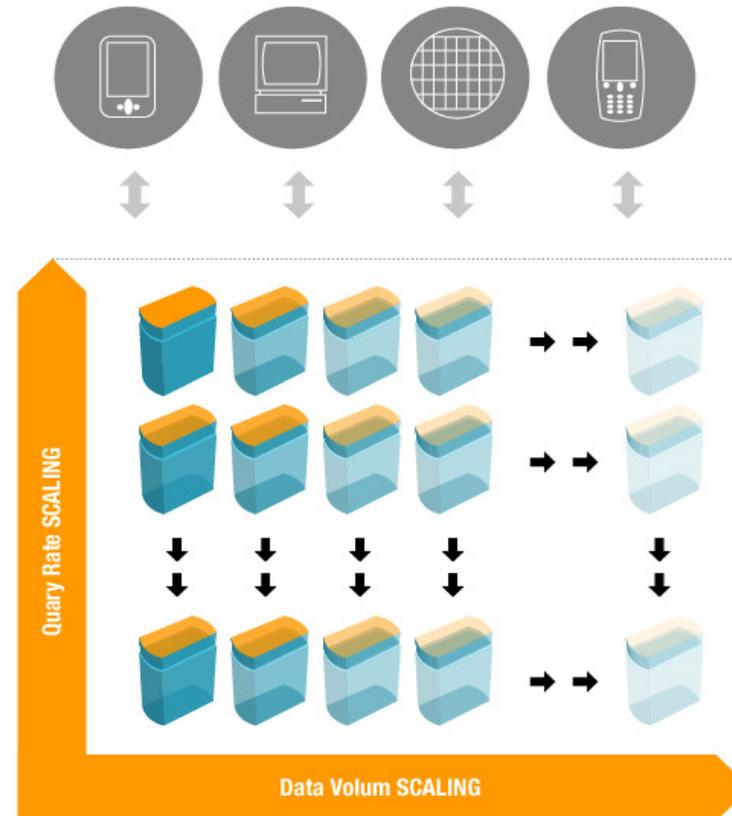
Relevancy

Completeness	How well does the query match superior contexts?
Authority	Is the document considered an authority for this query?
Statistics	How well does the content of this document match the query?
Quality	Is this a document of "high quality"?
Freshness	How old is this document, when was it last updated?
Geography	Where are you querying from?



Engine architecture

- **Scaling in data volume**
 - Add columns
 - Each column holds a partition of the data
 - Query the partitions in parallel
- **Scaling in query traffic**
 - Add rows
 - Replicate the partitions
 - Distribute the queries
- **Scaling in other dimensions**
 - Query complexity
 - Fault tolerance



Traditional Media Challenged

Entrance of Non-Traditional Competitors

Competitive Pressures Redefining the Business



Web Search Challenges

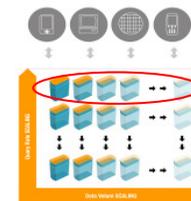
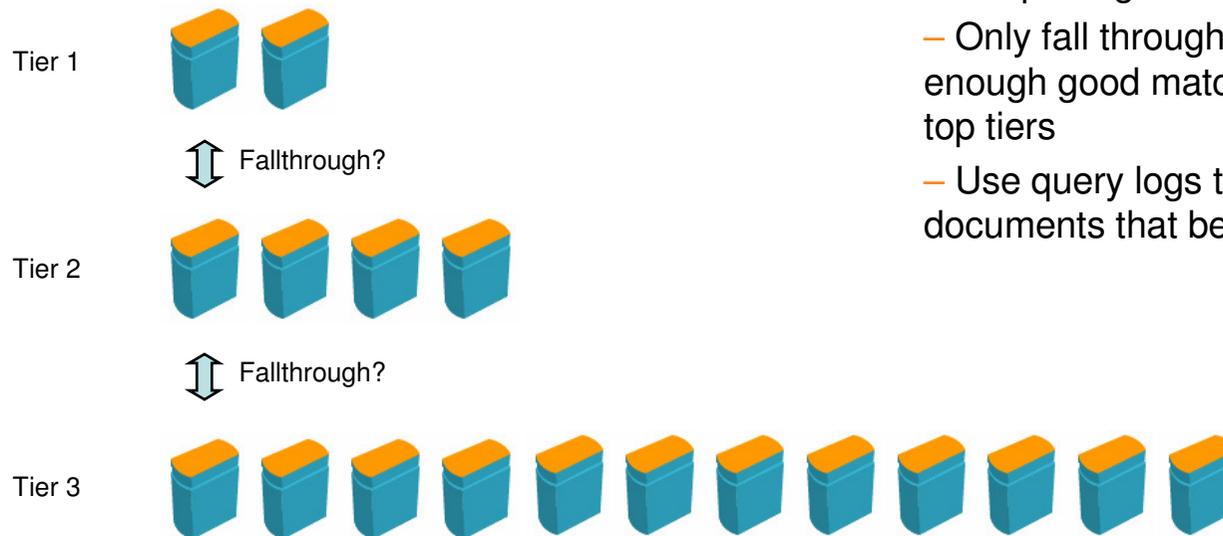
- **Business Model**
- **Architecture and operations**
 - Document volume, query volume and index freshness
- **Index quality**
 - Spam and offensive content detection
 - Duplicate detection
 - Algorithmic and/or editorial efforts
 - Characteristics must be defined
 - Filter or influence
- **Relevance**
 - Link cardinality
 - Anchor text
 - Freshness
 - Authorative sources
 - Editorial content

Engine architecture

- **For very big deployment scenarios**
 - Web scale, i.e., billions of documents

- **All search nodes are equal, but some are more equal than others**

- Organize the search nodes into multiple tiers
- Top tier nodes may have fewer documents and run on better hardware
- Keep the good stuff in the top tiers
- Only fall through to the lower tiers if not enough good matches are not found in the top tiers
- Use query logs to decide which documents that belong in which tiers



www.scirus.com



SCIRUS
for scientific information only

- **scirus.com – The Web’s Science Search**
- **Scientific classification**
- **Combining web and corporate index**

Search Within [Search Tips](#)

Information Types	Information Sources	
<input checked="" type="checkbox"/> All	<input checked="" type="checkbox"/> All Journal sources	<input checked="" type="checkbox"/> All Web sources
<input type="checkbox"/> Articles	<input type="checkbox"/> Beilstein on ChemWeb	<input type="checkbox"/> US Patent Office
<input type="checkbox"/> Scientist homepages	<input type="checkbox"/> MEDLINE on BioMedNet	<input type="checkbox"/> Neuroscion
<input type="checkbox"/> Patents	<input type="checkbox"/> ScienceDirect	<input type="checkbox"/> E-Print ArXiv
<input type="checkbox"/> Conferences	<input type="checkbox"/> IDEAL	<input type="checkbox"/> Other
<input type="checkbox"/> Abstracts	<input type="checkbox"/> BioMed Central	

Subject Areas

<input checked="" type="checkbox"/> All	<input type="checkbox"/> Environmental Sciences	<input type="checkbox"/> Medicine
<input type="checkbox"/> Agricultural and Biological Sciences	<input type="checkbox"/> Earth and Planetary Sciences	<input type="checkbox"/> Neuroscience
<input type="checkbox"/> Astronomy	<input type="checkbox"/> Law	<input type="checkbox"/> Pharmacology
<input type="checkbox"/> Biosciences	<input type="checkbox"/> Life Sciences	<input type="checkbox"/> Physics
<input type="checkbox"/> Chemistry and Chemical Engineering	<input type="checkbox"/> Languages and Linguistics	<input type="checkbox"/> Psychology
<input type="checkbox"/> Computer Science	<input type="checkbox"/> Materials Science	<input type="checkbox"/> Social and Behavioral Sciences
<input type="checkbox"/> Economics, Business and Management	<input type="checkbox"/> Mathematics	<input type="checkbox"/> Sociology
<input type="checkbox"/> Engineering, Energy and Technology		

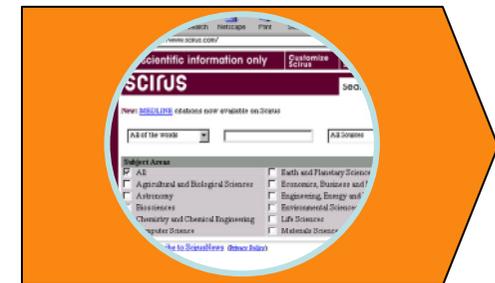
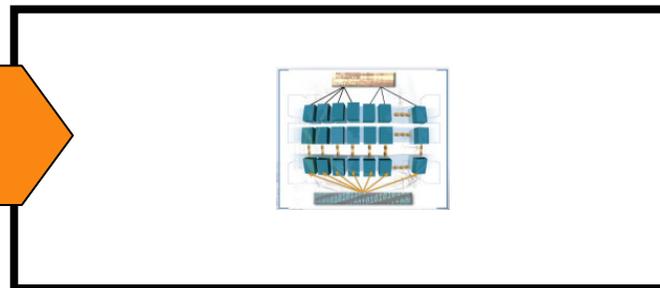
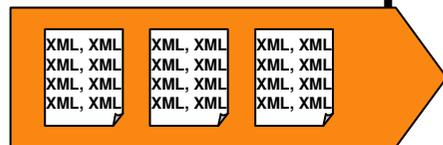
Results Published Between
1973 and 2002

Show 10 results per page

• You can save these settings as your search preferences, for use each time you visit Scirus.
• They can easily be edited or turned off at any stage.

[Search Tips](#)

- Refine your search using the following terms
- [cell injury](#)
 - [cerebral infarct](#)
 - [enlarged](#)
 - [etiology](#)
 - [health center](#)
 - [infarct](#)
 - [infarct of the myocardium](#)
 - [infarction](#)
 - [intestinal](#)
 - [myocardial](#)
 - [myocardial infarct](#)
 - [myocardium](#)
 - [necrosis](#)
 - [relational database](#)
 - [school of medicine](#)
 - [transverse](#)



90 M web pages
15 M Elsevier Science publications

- **Scientific classification**
- **Grouping and identification of related articles**

- **Leading science Index**
- **Understanding content**
- **Scientific navigation**

XML search

- **Moving beyond a flat document model**
 - From a simple (field, value) layout to complex, nested structure having scopes/tags
 - From a predefined index layout to schema flexibility
- **Some queries cannot be adequately handled without structure**
 - Flattening out the content won't quite work
 - False positives slip through

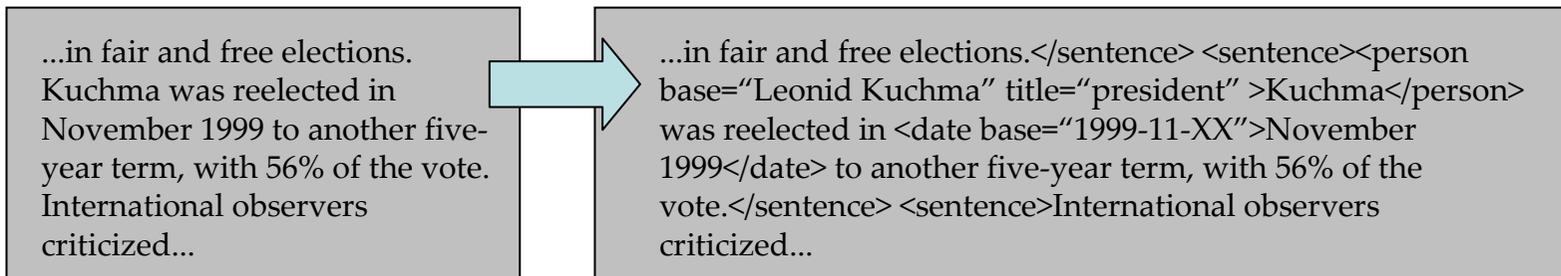
```
<authors>
  <author>John P. Brown</author>
  <author>George Smith</author>
</authors>
```

“Show me documents authored by John Smith”

```
- <PLAY>
  <MAINTITLE>The Tragedy of Antony and Cleopatra</MAINTITLE>
+ <FM>
+ <PERSONAE>
  <SCNDESCR>SCENE In several parts of the Roman empire.</SCNDESCR>
  <PLAYSUBT>ANTONY AND CLEOPATRA</PLAYSUBT>
- <ACT>
  <TITLE>ACT I</TITLE>
- <SCENE>
  <TITLE>SCENE I. Alexandria. A room in CLEOPATRA's palace.</TITLE>
  <STAGEDIR>Enter DEMETRIUS and PHILO</STAGEDIR>
+ <SPEECH>
+ <SPEECH>
+ <SPEECH>
+ <SPEECH>
  <STAGEDIR>Enter an Attendant</STAGEDIR>
+ <SPEECH>
+ <SPEECH>
- <SPEECH>
  <SPEAKER>CLEOPATRA</SPEAKER>
  <LINE>Nay, hear them, Antony:</LINE>
  <LINE>Fulvia perchance is angry; or, who knows</LINE>
  <LINE>If the scarce-bearded Caesar have not sent</LINE>
  <LINE>His powerful mandate to you, 'Do this, or this;</LINE>
  <LINE>Take in that kingdom, and enfranchise that;</LINE>
  <LINE>Perform 't, or else we damn thee.'</LINE>
</SPEECH>
- <SPEECH>
  <SPEAKER>MARK ANTONY</SPEAKER>
  <LINE>How, my love!</LINE>
</SPEECH>
+ <SPEECH>
  -----
```

Information extraction

- **Apply text mining techniques to identify entities of interest**
 - Structural and semantic regions
 - Makes unstructured data more structured
- **Mark them up in context**
 - Grammars as scope producers
 - Scopes can be annotated with meta data



- **Make it possible to act on the information!**
 - E.g., make it searchable in a way that preserves context

Scope

Scalable search

"Sentences where someone says something positive about Adidas."

```
xml:sentence:("adidas" and sentiment:@degree:>0)
```

"Paragraphs that contain quotations by George W. Bush, where he mentions a monetary amount."

```
xml:paragraph:quotation:(@speaker:"bush" and scope(price))
```

"Paragraphs that discuss a company merger or acquisition."

```
xml:paragraph:(string("merger", linguistics="on") and scope(company) and scope(price))
```

"Quotations where somebody says something about the Gaza Strip."

```
xml:quotation:("gaza strip")
```

XML technologies

Information extraction

"Sentences where the acronym 'MIT' is defined."

```
xml:sentence:acronym:(@base:"mit" and scope(@definition))
```

"Dates and locations related to the query 'd-day'."

```
xml:sentence:("d-day" and (scope(date) or scope(location)))
```

fast 

 enabling information on demand™



XML search

- **Queries that combine structure and content**
 - Impose contextual constraints on the content
- **FAST Query Language (FQL)**
 - Partial overlap with XQuery
 - Linguistic extensions
- **Return matching scopes**
 - See the matching document fragments, including markup and annotations

```
xml:sentence(("cancer" and scope(person))  
            ~  
            //sentence[fast:matches(., "cancer") and ./person])
```



```
<matches>  
  <match>  
    <sentence>The publicist of <person>Robert De  
Niro</person> announces that the actor has prostate  
cancer.</sentence>  
  </match>  
  <match>  
    <sentence><person>De Niro</person> was diagnosed  
with cancer last week.</sentence>  
  </match>  
  <match>  
    <sentence>"He'll fight the cancer," says  
<person>John Barnes</person>, founder of his Welsh  
fanclub.  
  </match>  
</matches>
```

Anatomy of a Search Platform

