

**MASSACHUSETTS INSTITUTE OF TECHNOLOGY
SLOAN SCHOOL OF MANAGEMENT**

15.565 Integrating Information Systems:

Technology, Strategy, and Organizational Factors

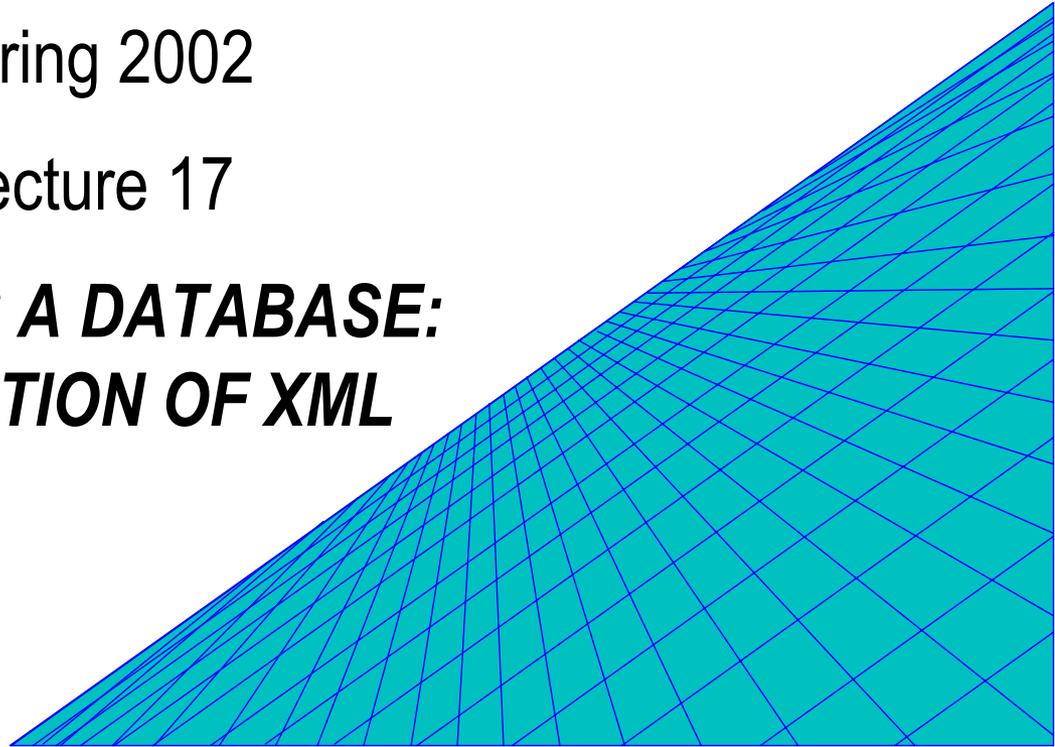
15.578 Global Information Systems:

Communications & Connectivity Among Information Systems

Spring 2002

Lecture 17

***WEB AS A DATABASE:
EVOLUTION OF XML***



Browser

←
“raw information”



Web sites

Traditional Use of the Web: For direct human usage Focus: **Entertainment**

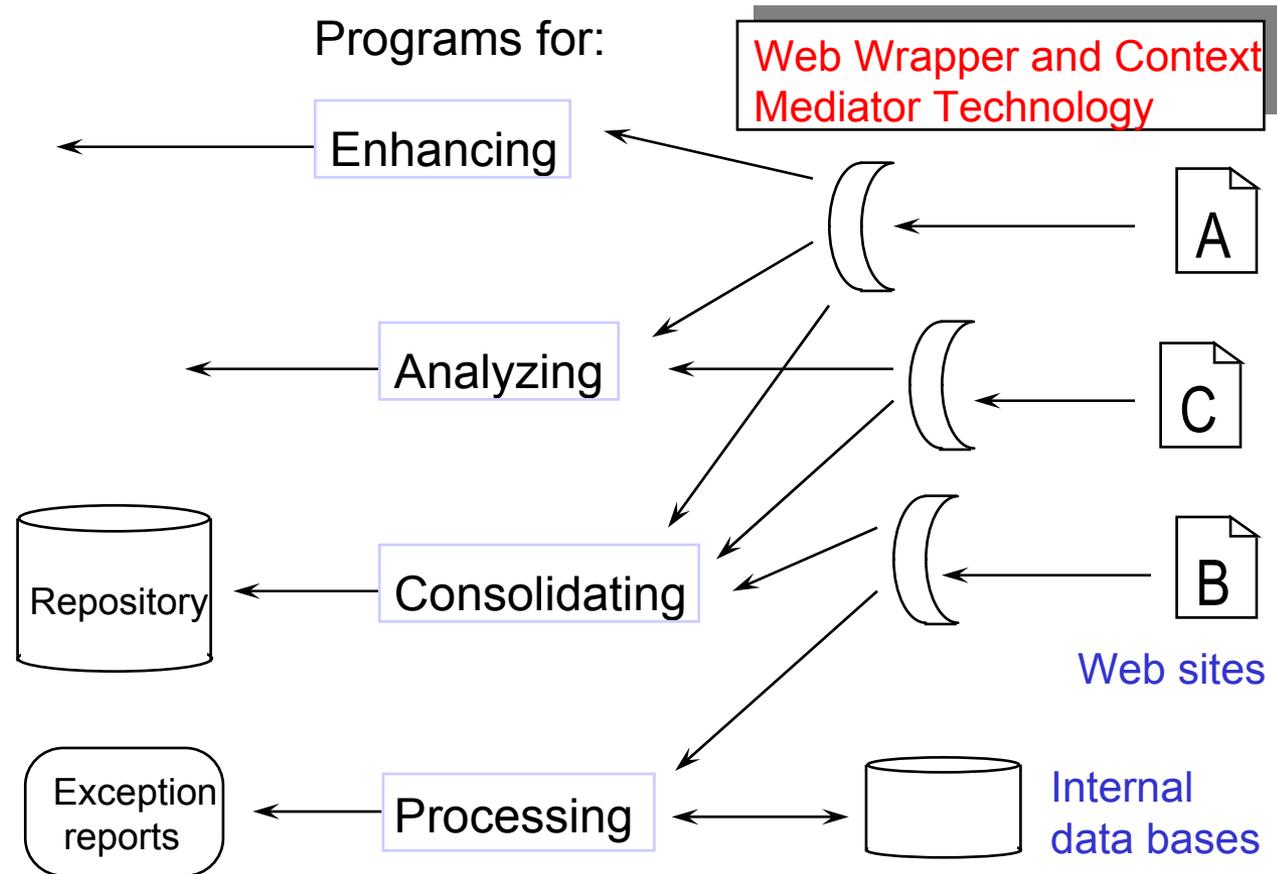
Examples:

Extract just_mortgage rate information

Compare mortgage rates offered by multiple sources

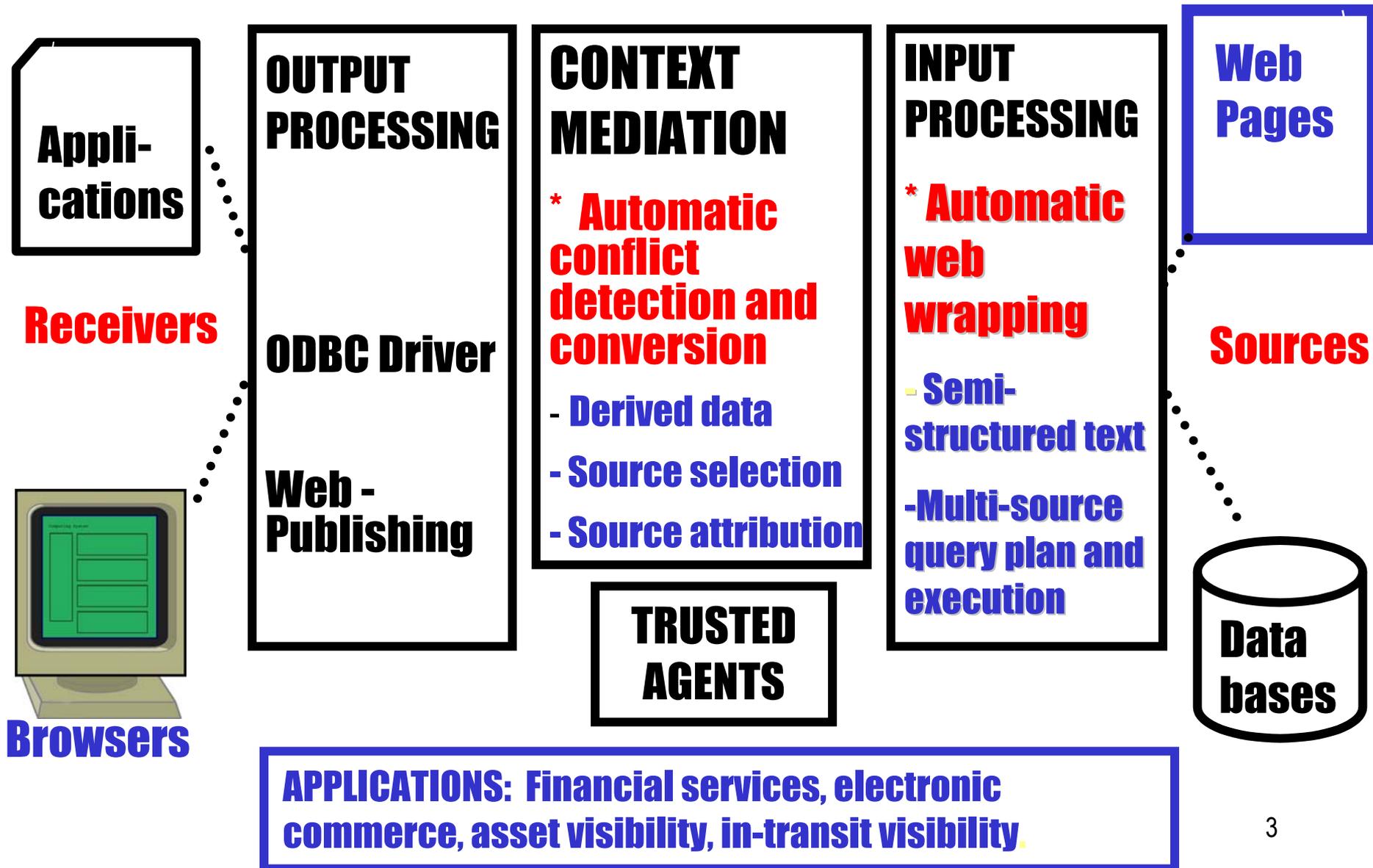
Build **cumulative** database of mortgage rates over time

Compare cumulative rates with previously stored, **alert** to new highs and lows



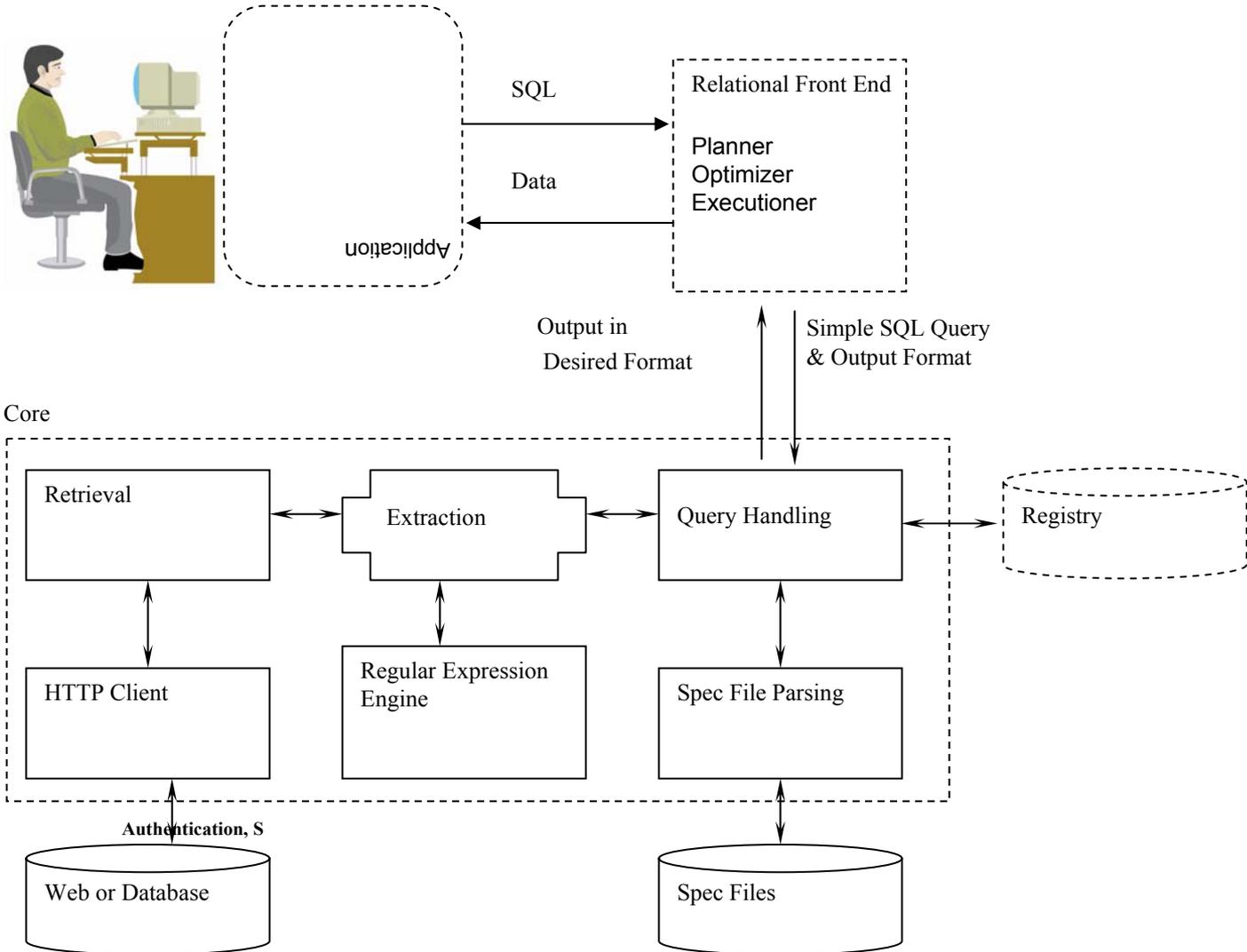
New Use of the Web: Program intermediaries Focus: **Productivity**

MIT Sloan COntext INterchange (COIN) Project



Example Semi-structured Web data: Intel SEC filing

Cameleon Architecture



CIA Fact Book Example

Singapore

Introduction

[\[Top of Page\]](#)

Background: Founded as a British trading colony in 1819, Singapore joined Malaysia in 1963, but withdrew two years later and became independent. It subsequently became one of the world's most prosperous countries, with strong international trading links (its port is one of the world's busiest) and with per capita GDP above that of the leading nations of Western Europe.

Geography

[\[Top of Page\]](#)

Location: Southeastern Asia, islands between Malaysia and Indonesia

Geographic coordinates: 1 22 N, 103 48 E

Map references: Southeast Asia

Area:

total: 647.5 sq km

land: 637.5 sq km

water: 10 sq km

Area - comparative: slightly more than 3.5 times the size of Washington, DC

Land boundaries: 0 km

Coastline: 193 km

Maritime claims:

exclusive fishing zone: within and beyond territorial sea, as defined in treaties and practice

territorial sea: 3 nm

Climate: tropical; hot, humid, rainy; no pronounced rainy or dry seasons; thunderstorms occur on 40% of all days (67% of days in April)

Terrain: lowland; gently undulating central plateau contains water catchment area and nature preserve

Elevation extremes:

lowest point: Singapore Strait 0 m

highest point: Bukit Timah 166 m

CAMELEON QUERY:

Select capital, location, coordinates, totalarea, climate, population, GDP
from cia where Country="Singapore"

CAMELEON RESULTS:

Record 1

CAPITAL

Singapore

LOCATION

Southeastern Asia, islands between Malaysia and Indonesia

COORDINATES

1 22 N, 103 48 E

TOTALAREA

647.5 sq km

CLIMATE

tropical; hot, humid, rainy; no pronounced rainy or dry seasons; thunderstorms occur on 40% of all days (67% of days in April)

POPULATION

4,151,264 (July 2000 est.)

GDP

\$98 billion (1999 est.)

Spec file for CIA Fact Book (partial)

#Relation=cia

#Source=http://www.odci.gov/cia/publications/factbook/country.html

#Attribute=Link#String

#Begin=Top\s*of\s*Page

#Pattern=#Country#

#End=</[Bb][oO][dD][yY]>

#Source=http://www.odci.gov/cia/publications/factbook/#Link#

#Attribute=Telephone#String

#Begin=Telephones:

#Pattern=\s*([\0-\377]*?)\s*<p>

#End=Telephone system:

#Attribute=Background#String

#Begin=Background:

#Pattern=\s*([\0-\377]*?)\s*<

#End=Location:

#Attribute=Location#String

#Begin=Location:

#Pattern=\s*([\0-\377]*?)\s*<p>

#End=Geographic\s*coordinates:

...

Regular Expressions Used in Spec Files

- * Match 0 or more times (greedy).
- *? Match 0 or more times (non-greedy).
- + Match 1 or more times (greedy).
- ? Match 0 or 1 time (greedy).

Greedy quantifiers such as * matches as much as possible, whereas non-greedy quantifiers stop at the minimum match. Example:

` hello <i>lovely </i> world `

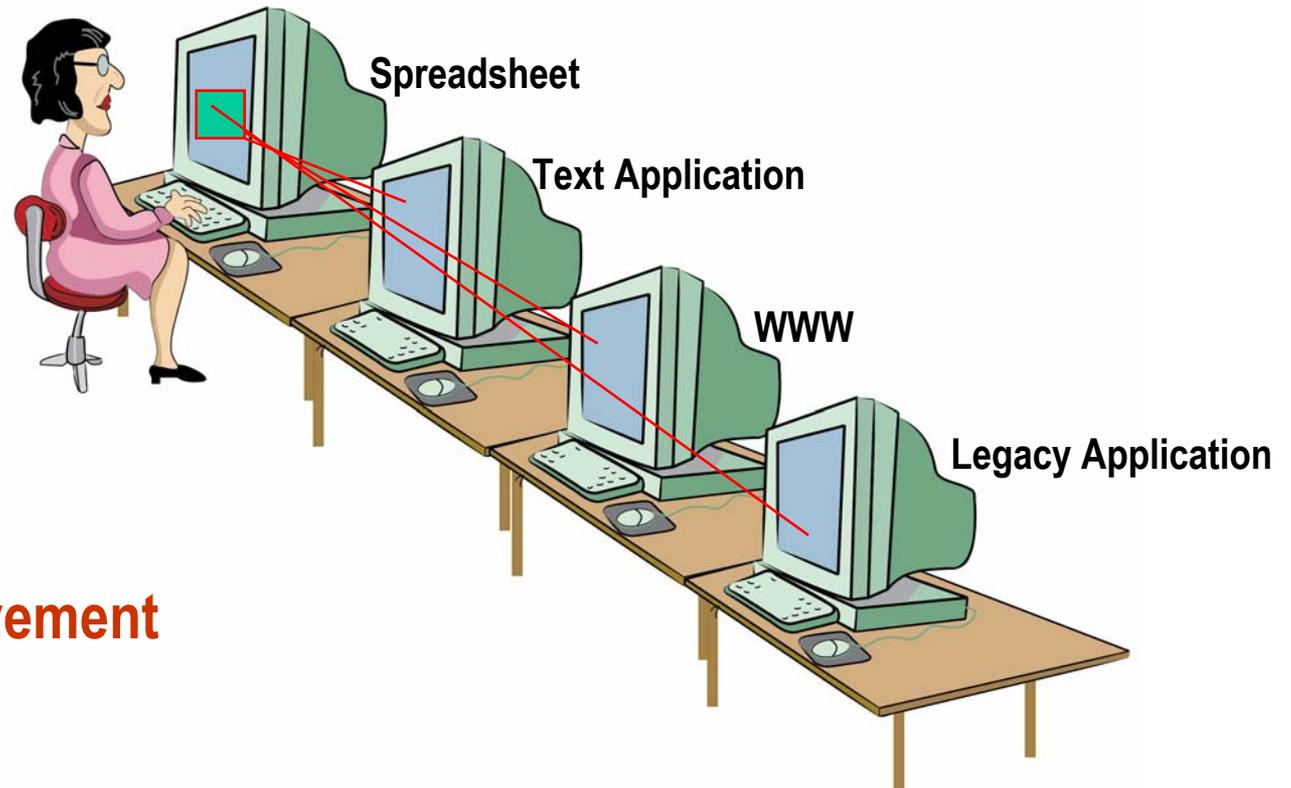
`(.*)` would match 'hello <i>lovely </i> world' whereas

`(.*?)` would match 'hello' and 'world'

- . matches everything except \n
- `[0-377]` matches everything
- ^ matches the beginning of a string or line
- `[^ a character]` matches everything except the specified character.
For instance `[^<]` matches anything but <
- \$ matches the end of a string or line
- \s matches a whitespace character
- \S matches a non-whitespace character
- \d matches a digit
- Expressions within parentheses are saved.

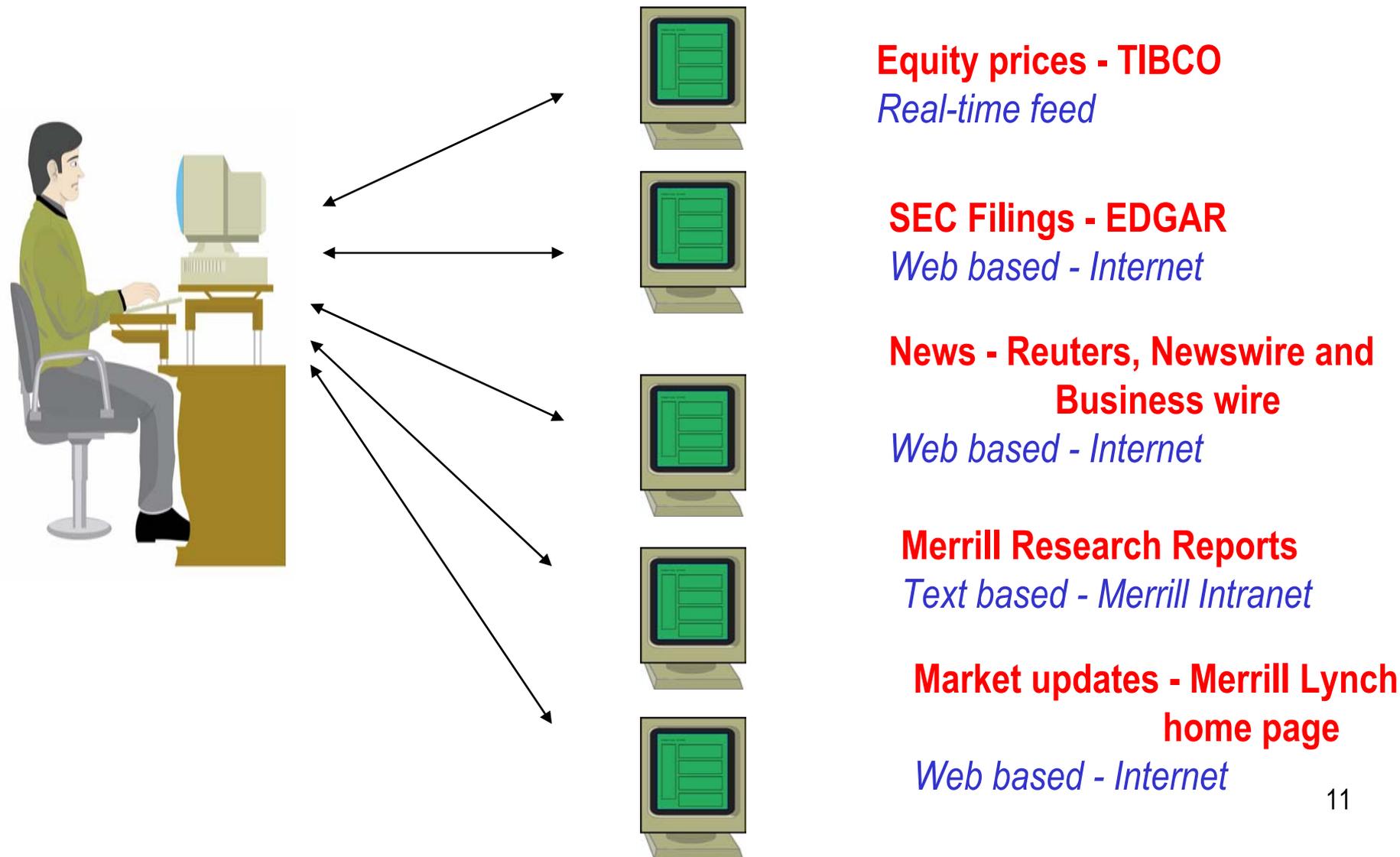
Sample Application

Research Analyst
or
Trader



Manual Data Movement

Providing Integrated Data and Analysis



Spreadsheet Interface



Real time
Tibco

Stock Portfolio		Attributes from Research Report			Attributes from Edgar SEC Filings				
1	Available Portfolios:								
3	Portfolio #2	Investment Opinion			Net Income				
5	Portfolio #3	52 - Week Range			Interest Expense				
6		Update			Clear				
8					<p>Market Value</p> <p>GM 36% INTC 22% IBM 41%</p>				
13									
19									
20	Stock	Ticker Symbol	Shares held	Purchase Price	Current Price	Market Value	Gain / Loss	Investment opinion	Net Income
21	Intel Crop	INTC	15000	50	93.625	1404375	654375	B-1-1-7	1983
22	Intl Bus Machine	IBM	24000	65	105	2520000	960000	B-3-3-8	1,195
23	General Motors	GM	35000	50	64	2240000	490000	B-2-2-7	1,796
24									
25					Total Value	6164375			
26	Latest News:								
27	Intel Crop	PRESS DIGEST - Wall Street Journal - July 29 - 2:12 am							
28	Intl Bus Machine	Eloquent to offer presentation system for Internet - Sunday November 2, 1997 - 12:32pm							
29	General Motors	Suzuki To Appeal Jury Award - Saturday November 1, 1997 - 12:13am							
30									
31									
32									

Free EDGAR™ COMPANIES [FILINGS] WATCHLIST ANALYSIS REFERENCE CONTACT HELP

Start a new Filings Search

Filings for INTEL CORP

View Company Profile for INTEL CORP

Today's SEC Filings

Content's...

Header:

Body:

10-9

Part 1

Item 1: Financial Statements

	215	76
Interest income and other, net	---	---
Income before taxes	3,075	1,376
Provision for taxes	1,092	482
Net income	\$ 1,983	\$ 894
Earnings per common and common equivalent share	\$ 2.20	\$ 1.02
Cash dividends declared per common share	\$ 0.05	\$ 0.04
Weighted average common and common equivalent shares outstanding	900	880

Home Off Search Help Menu

Get 20% cash rebates at 1000s of restaurants nationwide. 30-day free Holiday to Card membership. Click here.

[Business][Financial][Company][Industry][PR Newswire][Business Wire][Quotes]

Search News Help

Intel Corp

Company Profile - Current Stock Price

Tuesday July 29, 1997

PRESS DIGEST - Wall Street Journal - July 29 - Reuters - 2:12 am

Merrill Lynch Research Report: General Motors

Long Term Recommendation: ACCUMULATE

Price: \$57

12 Month Price Objectives: \$68

	1996A	1997E	1998E
Estimates (Dec)			
EPS:	\$7.44	\$7.57	\$8.23
P/E:		7.63x	7.87x
Opinion & Financial Data			
Investment Opinion			B-2-2-7
Mkt. Value / Shares Outstanding (mn) :		\$89,011	

XML – The Silver Bullet ?

- XML is (according to press reports ...)
 - “HTML on steroids” ?
 - “a Rosetta Stone” ?
 - “a universal way to translate data” ?
 - “a miraculous way to” ... information integration ?
 - “a silver bullet” ?

XML What is it?

- **XML** - **EX**tensible **M**arkup **L**anguage
- Meta language for defining a markup language
- Based on SGML - Standard Generalized Markup Language
- Data model for syntax for structuring data
- Can define tags at will
- Can nest document structures to arbitrary levels of complexity
- Can use Document Type Definition (DTD)
- Many other members of “family”:
 - XSL, XSLT, XLL, XML-Query, etc.

XML Does help create structured Web pages

<u>Feature</u>	<u>HTML</u>	<u>XML</u>
Extensibility	Fixed set of tags	Extensible set of tags
Tag purpose	Presentation	Content
Views	Single	Multiple (XSL)
Orientation	Documents	Documents + semi-structured data
Search	Keyword	Keyword plus Field-sensitive queries

Example: HTML Compared with XML

HTML *

```
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 3.2 Final//EN">
<html>
<head> . . .
<BODY topmargin=18 leftmargin=6 bgcolor="#ffffff" link="#0000ee" VLINK="#551A8B" ALINK="#ff0000">
<pre><font size=2>
                Regular   Our
                Price     Price
    Palm Pilot V  329.00   236.00   In stock
</font></pre>
<table cellpadding=0 cellspacing=0 border=0>
<tr><td align=left valign=middle width=455 nowrap height=20>
<tr><td align=left valign=top nowrap width=455>
<font size=1 face="helvetica,arial"> . . .
</BODY>
</HTML>
```

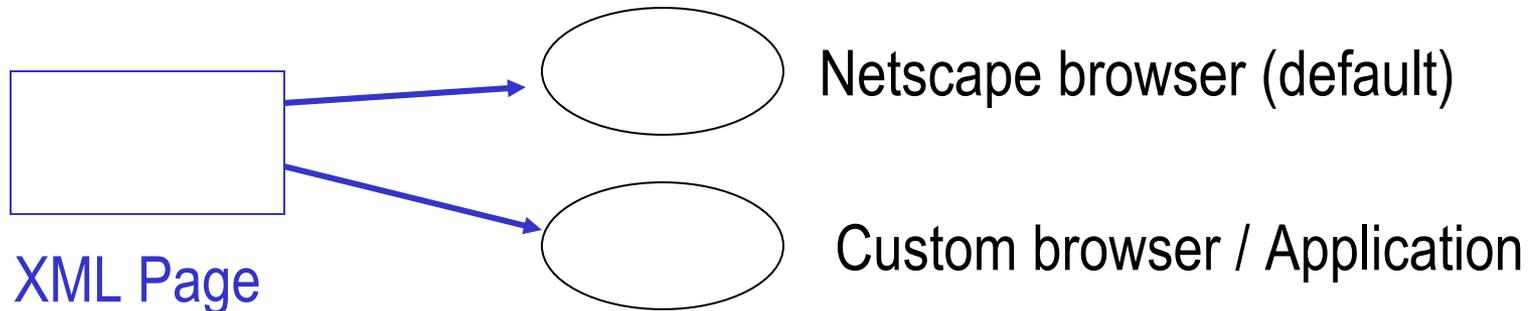
* The HTML is normally much messier, with lots more formatting details and `<tr>` and `<td>` tags for defining tables and tab positions.

XML

```
<XML>
<Product info>
  <Product> Palm Pilot V </Product>
  <Regular price> 329.00 </Regular price>
  <Our price> 236.00 </Our price>
  <InStock> yes </InStock>
</Product info>
</XML>
```

XML Why do we need it?

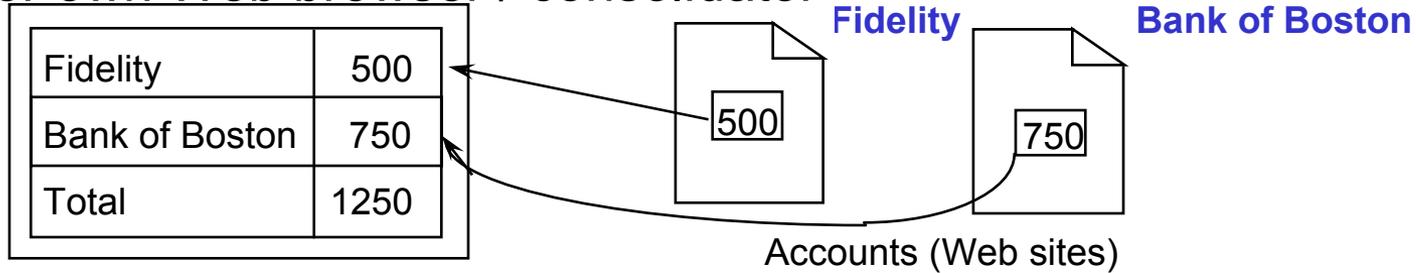
- W3C wanted to get out of the tag creation business
- Separate data from presentation
 - Use of a style sheet instead of “hardcoded” HTML formatting
 - Flexibility / scalability / extensibility



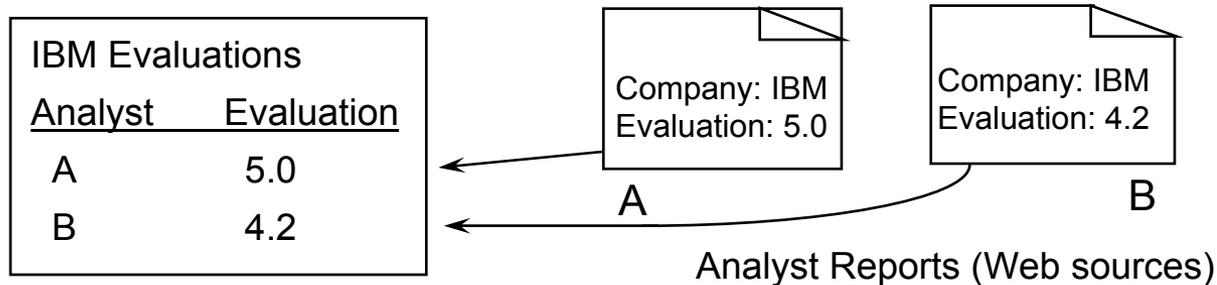
- Also important for Wireless Applications (WML/ XHTML)
- Human readability
- Computer processable
- Information interchange

Sample Applications of Semi-Structured Web Data

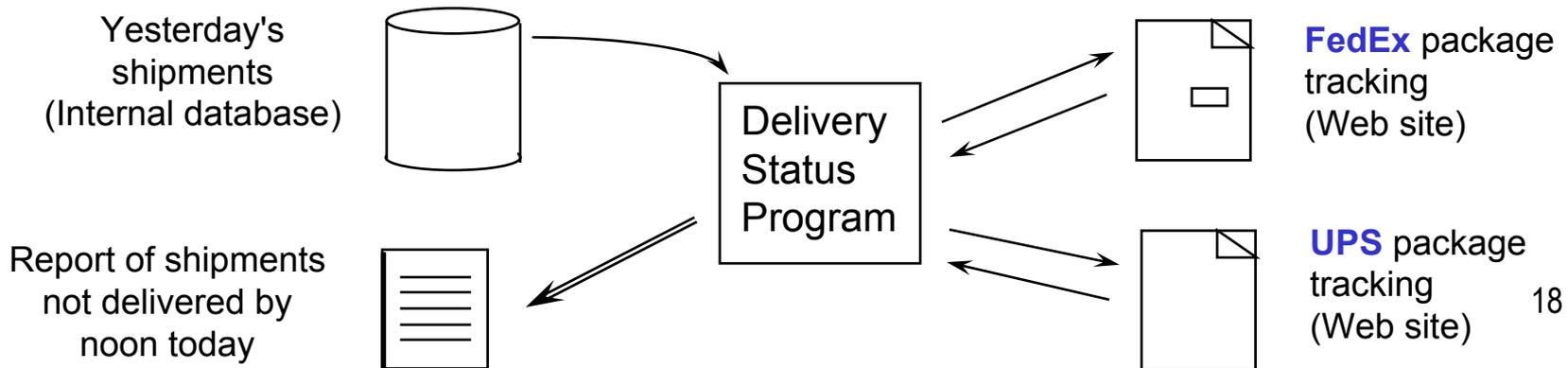
- Automate **Extraction** of data from specific Web sites into user tool, like Excel, or own Web browser / consolidator



- Automatically **Select** and **Consolidate** information across Web sites



- Integrate** Internet / Intranet / Client Server networks for internal operations



XML . . . Multiple Standards

- What is so great about XML standards – is that there are so many . . .
- Should tag for catalog be called “price” or “cost” ?
- “The director of electronic trading at Credit Suisse First Boston and chairman of financial services XML working group is struggling with [more than a dozen XML protocols . . . for financial trading applications.](#)” (*ComputerWorld*, 9 July 2001)

XML – The Silver Bullet ?

- XML is not quite:
 - “a Rosetta Stone”
 - “a universal way to translate data”
 - “a miraculous way to” ... information integration
 - “a silver bullet”
- It is a helpful tool toward information integration . . .
- Some background sources: w3c.org/XML and XML.org
- But much more is needed for information integration
 - Context Interchange and Semantic Web research are promising areas . . .

Summary

- **Tim Berners-Lee**, W3C Director:
 - "The Web is quickly becoming the world's fastest growing repository of data"
- In past: Primarily processed by humans.
- In future, must be processible by programs (agents for humans)
- Tools, such as MIT's Automatic Web Wrapper and W3C's XML, are providing these capabilities.

Work reported herein has been supported, in part, by the USA Advanced Research Projects Agency, Banco Santander Central Hispano, Citibank, Fleet Bank, First Logic, Merrill Lynch, PricewaterhouseCoopers, MIT Total Data Quality Management Program, MIT Center for eBusiness, Suruga Bank, and USAF/Rome Laboratory.