# Business Process Data Warehousing: Modeling and Integration Issues

Umeshwar Dayal
HP Fellow & Co-Director,
Advanced Business Intelligence Lab, HP Labs
umeshwar.dayal@hp.com

Fabio Casati, Malu Castellanos,
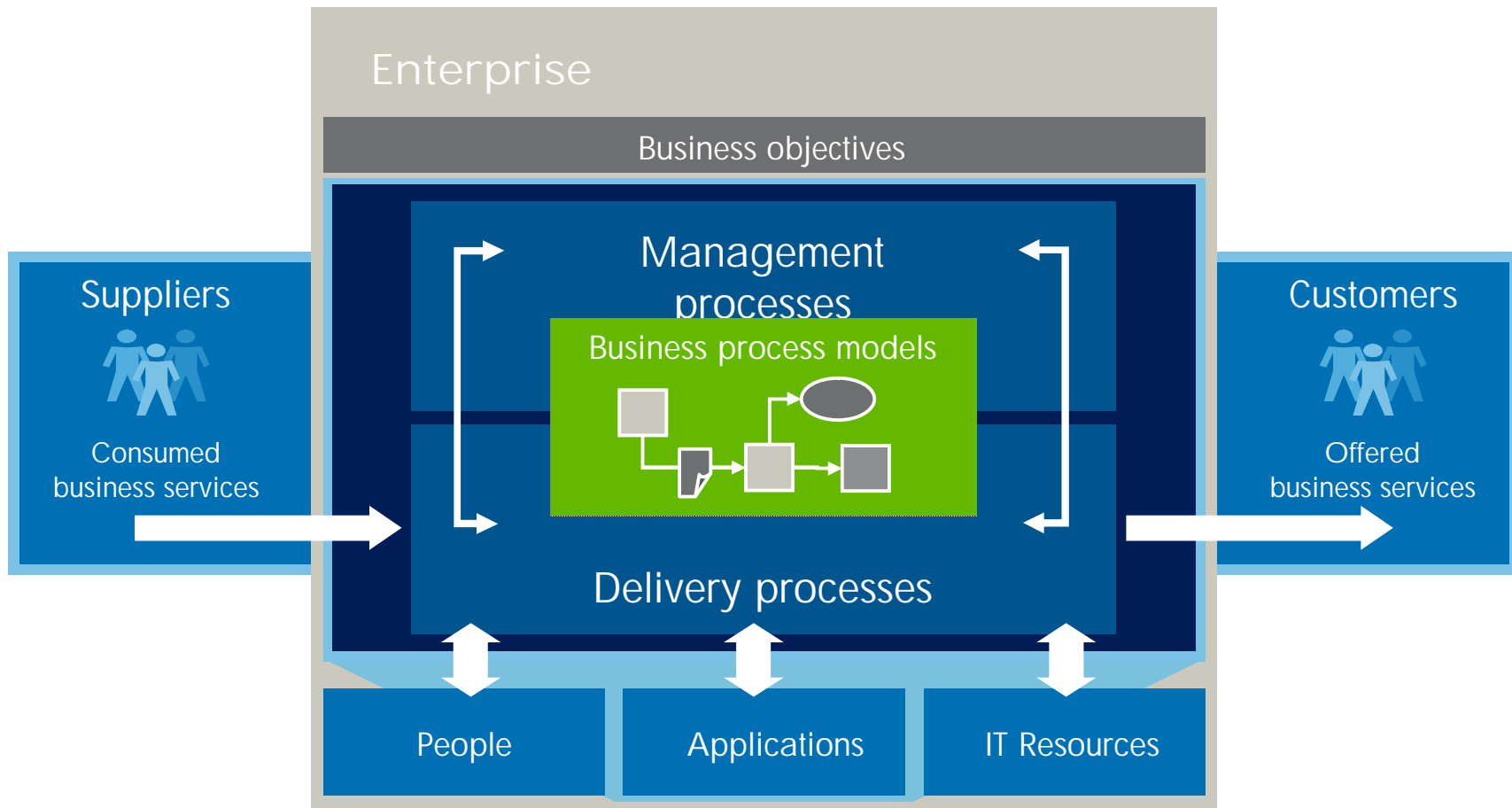Ming-Chien Shan, Norman Salazar,
Mehmet Sayal

INFINT Workshop, Bertinoro
September 2007

# Outline

- Context
  - Business Process Management
  - Business Process Intelligence
  - Relevance of Information Integration
- Process Modeling Issues
  - Process Views
  - Metrics Model
- Information Integration Issues
  - Generic Data Warehouse Schema
  - Abstraction Mechanisms
  - Generic ETL
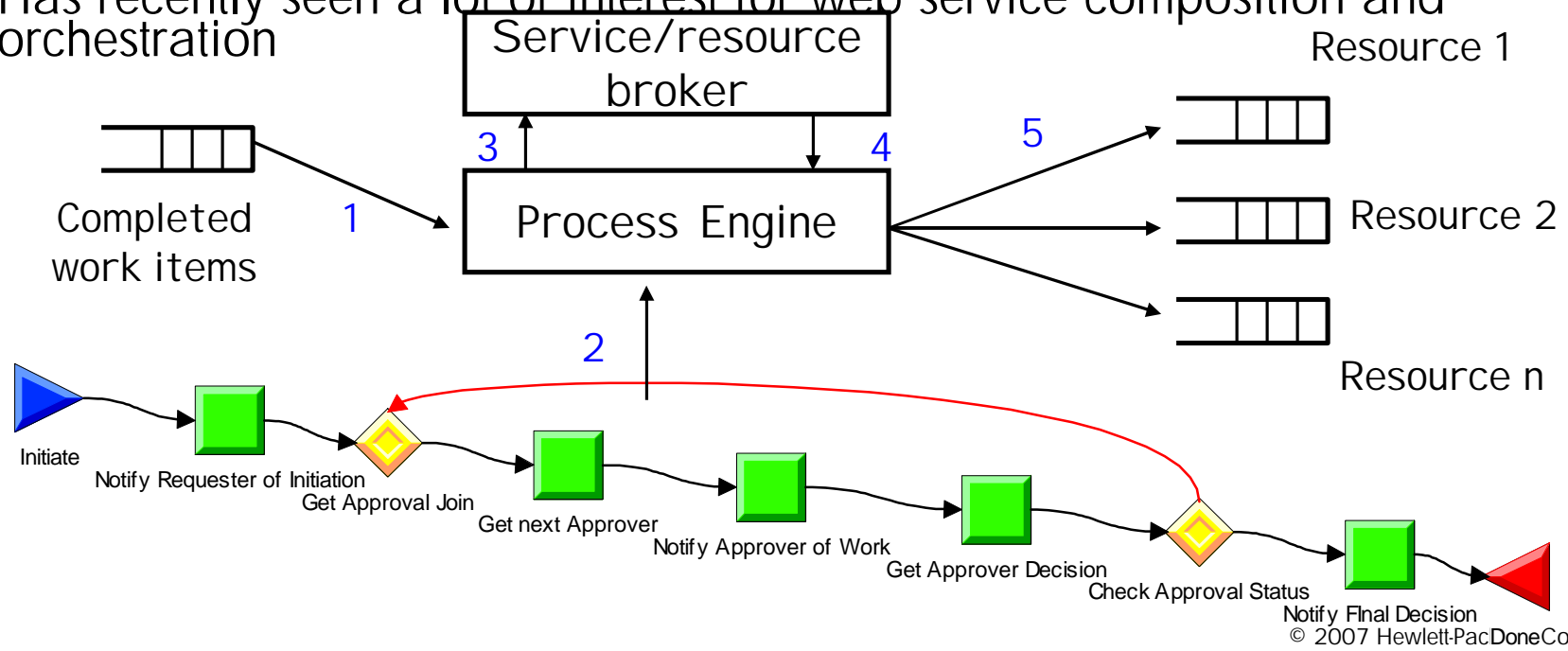  - Information Extraction from Semi-Structured Data
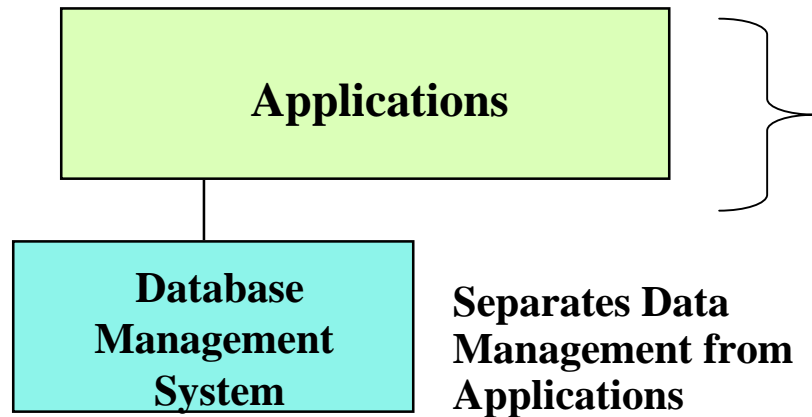- Summary

# Business Processes Drive the Enterprise



Supply chain processes: Procurement, Inventory, Logistics
Manufacturing processes: ERP, Product Lifecycle
Administrative processes: HR, Finance, Legal
Customer facing processes: support, CRM
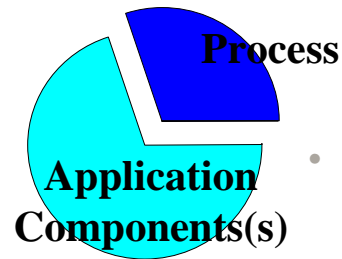
# Business Process Management

- Support the definition, execution, and monitoring of operational business processes.
- Define and control the sequence, timing, and other dependencies among tasks, as well as enforce business rules and policies.
- Assign resources (humans, data, applications, services) to individual tasks.
- Monitor and track all aspects of the process execution for auditing, business reporting, and process optimization.
- Has recently seen a lot of interest for web service composition and orchestration



Service/resource broker

Resource 1

Completed work items

3          4          5

1          Process Engine          Resource 2

2          Resource n

Initiate

Notify Requester of Initiation

Get Approval Join

Get next Approver

Notify Approver of Work

Get Approver Decision

Check Approval Status

Notify Final Decision

# Business Process Management Rationale

**Applications**

**Database Management System**

Separates Data Management from Applications

**Process**

**Application Components(s)**

- Process logic is "hard wired"
- Need programmer to change process

**Applications**

**Database Management System**

**Process Management System**

§ **Extract Rules for Process Flow and Resource Allocation from Applications and Data**

§ **Flexible and modular development to minimize impact of changes in one area on other areas:**

§ **Business Processes and Process Management become a corporate asset**

  § **a foundation for corporate business process reengineering and optimization**

# Example of Business Process Definition

**Main Process**



**Nested Sub process**



**Process definition captures:**

- Tasks (activities): manual, automated
- Control Flow
- Data Flow
- Exception handling
- Escalation
- Event triggers
- Resource resolution
- Application integration

# Resurgence of Interest in BPM

- First wave (1980's - 1990's) was focused on Business Process Automation/ Execution using Workflow Management technology: this had little success
  - Heterogeneous components
  - Cost
  - Complexity of building the WF system
  - Lack of support for application development lifecycle
  - Lack of standardization

- Second wave (2000's) is focused on Business Process Improvement: now there are the conditions for BPM to succeed
  - Web services and Service-Oriented Architectures
  - Maturity in the basic middleware, Open source BPMs
  - Standardization (BPEL, RosettaNet)
  - Trend towards increased Outsourcing
    - Understand and model the customer's and provider's business processes
    - Define and formalize Service Level Agreements (SLAs)
    - Alert/Predict SLA violations
    - Audit: you are liable, need to track
    - Analysis/optimization: much larger emphasis on operational efficiency
  - Regulatory compliance

# Business Process Intelligence

- Goal: improve the quality and performance of intra- and inter-enterprise business processes
  - Internal quality, as perceived by the service provider (e.g., reduced operating costs, fewer exceptions)
  - External quality, as perceived by the service consumer (e.g., better quality of service, reduced cost of service)
- Approach: Apply business intelligence techniques (data warehousing, data mining, simulation and optimization) to data relevant to business process execution
  - Integrate data from many sources:
    - Process management system (workflow engine) logs
    - web service execution logs
    - application logs
    - audit logs
    - event logs
    - systems management data, resource utilization data,
    - financial data, other business and operational data, …
  - Use the data to analyze, understand, and optimize processes
    - Resource assignments
    - Reporting on performance and quality of resources, service providers
    - Load prediction and optimization
    - Exception understanding and prevention
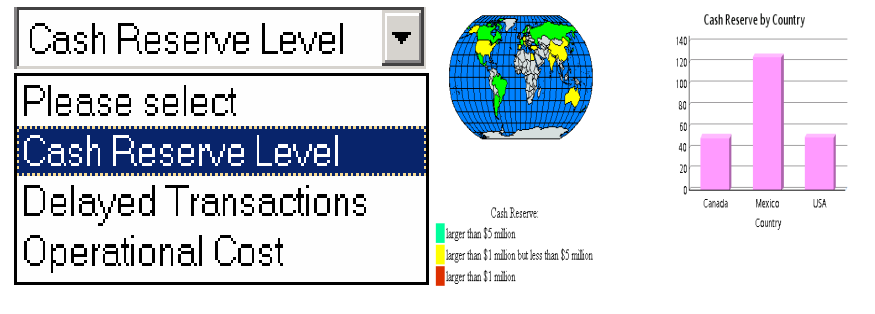    - Paths followed in the process graph

# Users want a crystal ball….

- How can I identify bottlenecks?
- What are the causes of missed SLAs (or other business performance metrics)?
- Can I predict my risk of missing SLAs (e.g., late payments)
- How much money do I save on electronic invoices versus paper ones? And how much time?
- How do the different payment methods compare in terms of cost and time?
- Can I predict my workload?
- What's my optimal resource plan? How many resources do I need to meet my SLAs (e.g., payment schedules)
- What's the disruption caused by unavailable resources?
- How do I "improve" my business process?
- What is the business impact of changing my IT infrastructure? My business process?

- Today this is difficult to do, requiring lots of custom design, system integration, and implementation effort.

# Business Process Intelligence
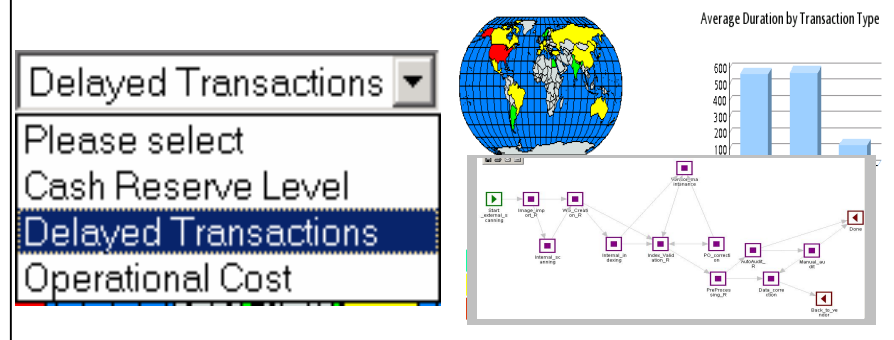# Embedding BI into Business Operations

## Financial Controller



Today, it is extremely difficult to automatically discover these inter-relationships and their implications !

## Operations Manager



## IT System Administrator

| Name | R-Eddie |
|---|---|
| Type | IT resource |
| Status | Critical |
| Change time | 2004-01-26 10:57:02.687 |
| Duration | 29 minutes 53 seconds |

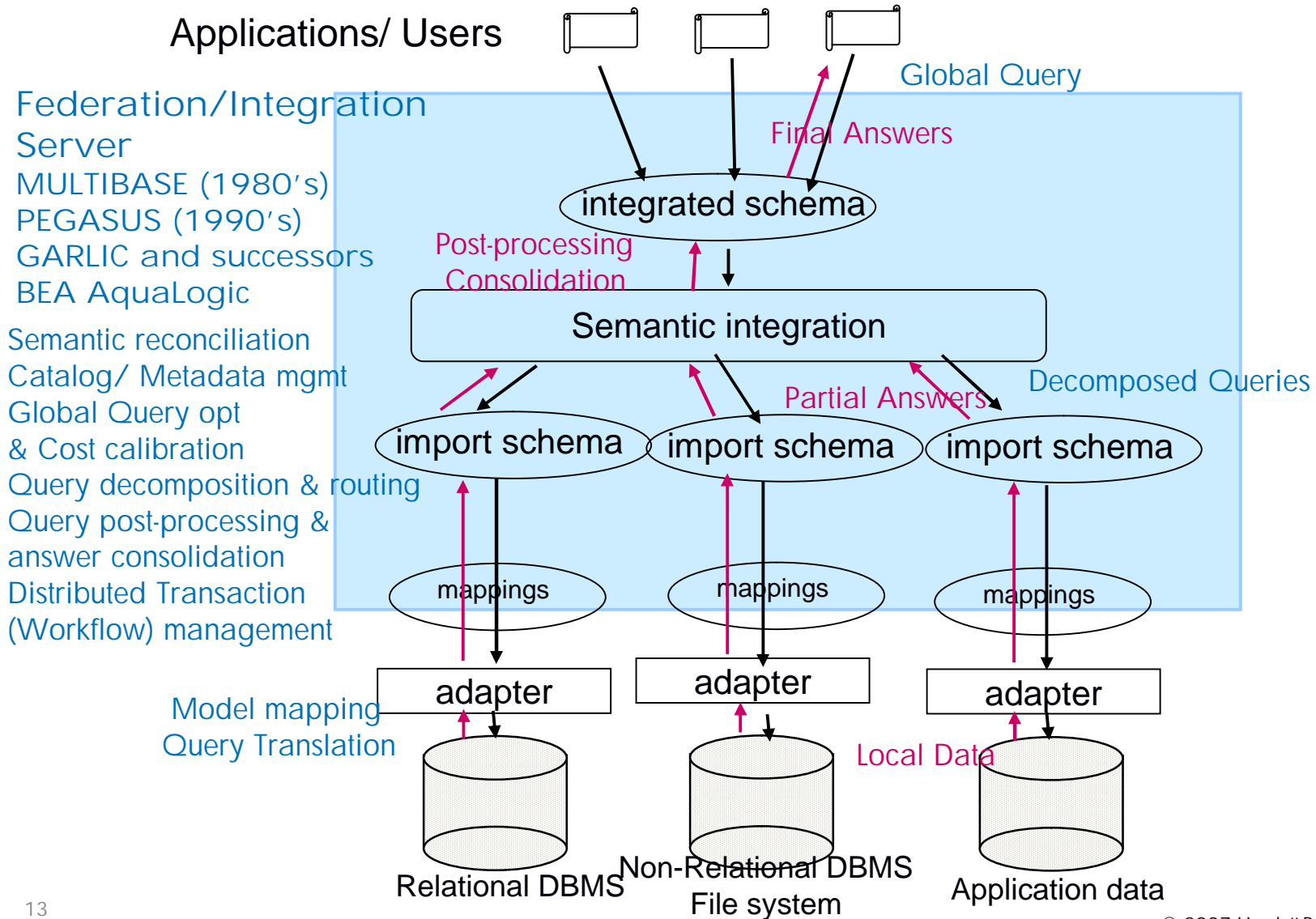# Suddenly, a very crowded space

Buzzwords popping up all over the place

- Business Service Management
- Business Event Management
- Business Activity Monitoring
- Business Operations Management
- Business Performance Management
- Business Process Intelligence
- Operational Business Intelligence
- Real-time enterprise, Zero latency enterprise
- Executive dashboards

# Information Integration Approaches

- Federated Databases (virtualization, stateless)
- Data Warehouses (materialized, stateful)
- Operational Data Stores
- Master Data Management
- Active Data Warehouses (hybrid of DW and ODS)

- For Business Process Intelligence, building a data warehouse (actually an active DW) is more appropriate
  - Need data from many sources (many of which are not databases)
  - Need historical data in addition to data about current process instances
  - Need complex transformations (e.g., map system events into abstract process progression)
  - Many reporting and analytic tools already work with data warehouses
- We built a research prototype, Business Cockpit, and tested it in several internal pilot solutions
- Several research challenges

# Federated databases: Data Virtualization

Applications/ Users

Federation/Integration Server
MULTIBASE (1980's)
PEGASUS (1990's)
GARLIC and successors
BEA AquaLogic

Semantic reconciliation
Catalog/ Metadata mgmt
Global Query opt
& Cost calibration
Query decomposition & routing
Query post-processing &
answer consolidation
Distributed Transaction
(Workflow) management

Model mapping
Query Translation

Global Query

Final Answers

integrated schema

Post-processing
Consolidation

Semantic integration

Partial Answers

Decomposed Queries

import schema    import schema    import schema

mappings    mappings    mappings

adapter    adapter    adapter

Local Data

Relational DBMS    Non-Relational DBMS
File system    Application data
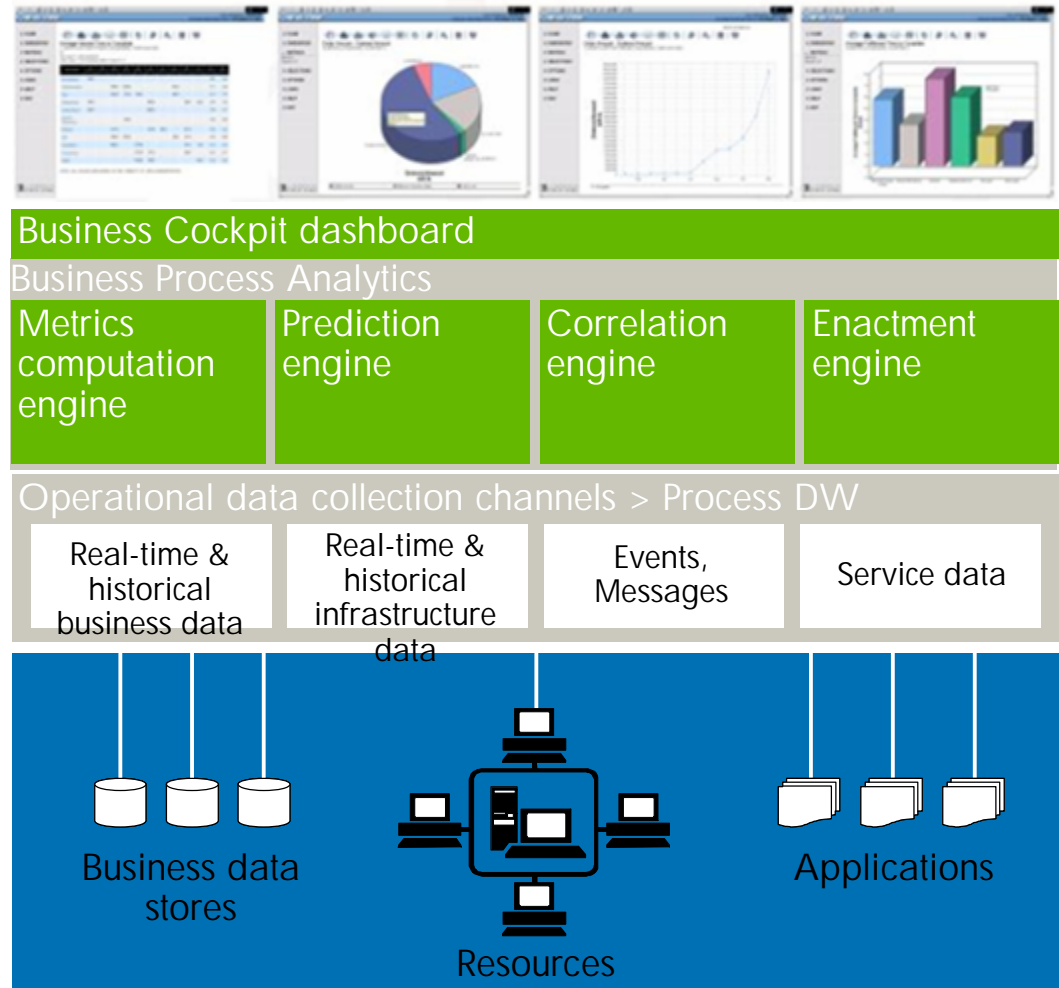
# Data Warehouse: Data Materialization
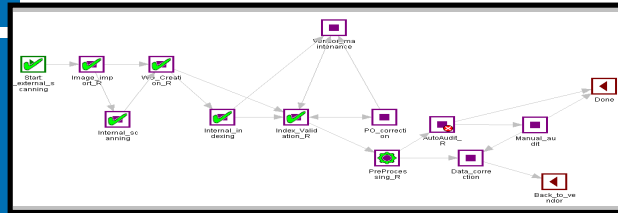


Source: IDC, 2005

# Business Cockpit: A Prototype Business Process Intelligence Solution

- Automatically correlate real-time and historical business process and IT systems data

- Customize/monitor/ analyze user-defined metrics

- Identify critical factors influencing business metrics (e.g., SLA compliance)

- Assess and predict impact of events on business KPIs

- Allocate IT capacity to optimize business operations



**Business Cockpit dashboard**

**Business Process Analytics**

| Metrics computation engine | Prediction engine | Correlation engine | Enactment engine |
|---|---|---|---|

**Operational data collection channels > Process DW**

| Real-time & historical business data | Real-time & historical infrastructure data | Events, Messages | Service data |
|---|---|---|---|

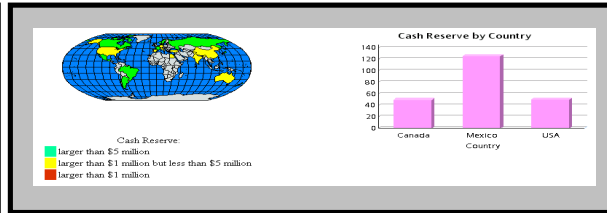Business data stores

Resources

Applications

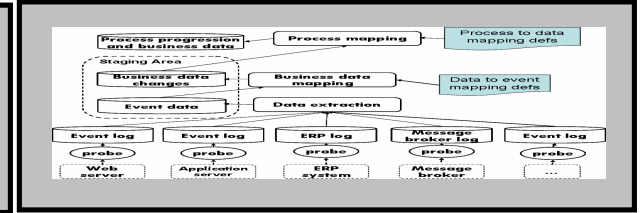# Business Process Intelligence Research Areas



**Process modeling**
Different views of the same process
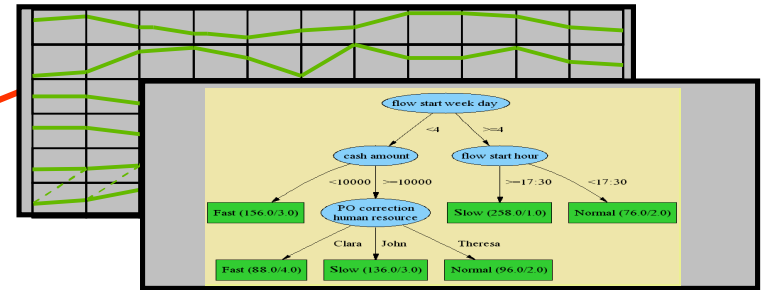Mappings between process views

**Metrics**
Functions/Queries over process and related data
Mappings between views

**Process Warehouse Design**
Generic and customizable schema and ETL suitable for any processes and data sources

## Business Cockpit

**Process Discovery**
Mine event logs to learn "true" implemented process
Automatically suggest process restructuring

**Simulation and Optimization**
Capacity management
Service/Resource selection

**Prediction**
Time-series prediction

**Correlation & Explanation**
Data mining to learn relationships between metrics, and between metrics and events

timeToNotifyAcceptance SLA is likely to be violated

16

# Outline

# Process Modeling: Abstract versus Implemented Process Views

Abstract process



Implemented process



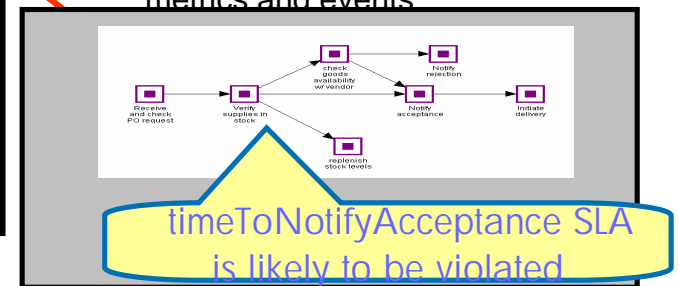- How to map between these views? Which is the base model and which is the view?
- Abstract process models are usually constructed manually
- No good tools for model refinement and implementation; the implemented process is often not explicitly modeled
- Process discovery: data mining techniques to learn and validate the implemented process
- Construct mappings between abstract and implemented process views (schema level, instance level, data level)
- Process integration and views over multiple processes are wide open problems

# Process Discovery and Validation



Discovery and validation

process logs

Web svc logs

app logs

resource logs

correlation

Monitoring tools

Business processes (composite services)

Web Services

Shared application infrastructure

virtualized resource pool
programmable, heterogeneous, distributed, shared

Management and control
(Web Services and business processes)

# Log Abstraction

Trace: a set of time ordered log entries that refer to the same process execution

AsBsAeCsBeAsAeCe …

BsBeAsDsAeDs …

CsCeAsAeEsFsFeEe …

- Log: a set of traces
  - Log supported by transaction monitoring systems
- Problems: missing data, wrong timestamps, incomplete traces, correlation among log entries

# Discovery of Ordering Constraints

- Possible process flow structures
  - Sequences
  - Splits
    - XOR
    - AND
    - OR
  - Joins
    - XOR
    - AND
    - OR
  - Loops

# Mining frequent sequences (episodes) from logs

Sequence

$(A) \rightarrow (B)$

AsAeBsBe

XOR-Split

$(A)$ → $(B)$
$(A)$ → $(C)$

AsAeBsBe

AsAeCsCe

AND-Split

$(A)$ → $(B)$
$(A)$ → $(C)$

AsAeBsCsBeCe
AsAeBsCsCeBe
AsAeCsBsBeCe
AsAeCsBsCeBe

# Model Discovery

- Progressive elimination
  - Probabilistic, iterative approach
  - Look at variance, average interval among activities
  - Look across traces
- It's not just the flow
  - lots of thresholds, tuned through classification
- Tuning thresholds is difficult

# Monitoring Processes: Types of Metrics

- Process metrics
  - execution times, durations, volumes, paths taken, outcomes
  - correlation with "previous" step
- Resource metrics
  - Performance of human and system resources in executing steps.
  - Correlation between resource and process metrics: which resources statistically lead to successful or unsuccessful executions, or which resources have led to certain paths being taken, e,g., escalations or error handling
- Business metrics
  - Domain-specific metrics, e.g., order-to-cash, turn around time, cash reserve levels
  - Correlation of business data with process data, e.g., efficiency and quality of execution based on invoice type.
  - Correlation between business data and resources, e.g., number of invoices from a given center processed by a given employee

# Conventional approach to defining and computing metrics

average **execution quality** by process

total **execution quality** by day

total **'time to acceptance' SLA violations** by customer

details on total **'time to acceptance' SLA violations**

predictions and explanations on **SLA violations**



| code (C, Java, SQL, OLAP) | code | ... | ... | ... | ... | ... | ... | | | code |

Process execution data

# Problems

- Long development time
  - Want to define 100 reports? Write 100 queries, test, deploy
- Poor performance
  - Want to view the reports? Run 100 queries
  - Real time? Concurrent users?
- Poor functionality
  - Limited support for drill-downs
  - No support for polymorphism
- Not robust to changes
  - Want a little different perspective? Write code, test, deploy
- No support for analytic features
  - Root-cause analysis, predictions, impact,…

# Our Approach: Provide Modeling Abstractions

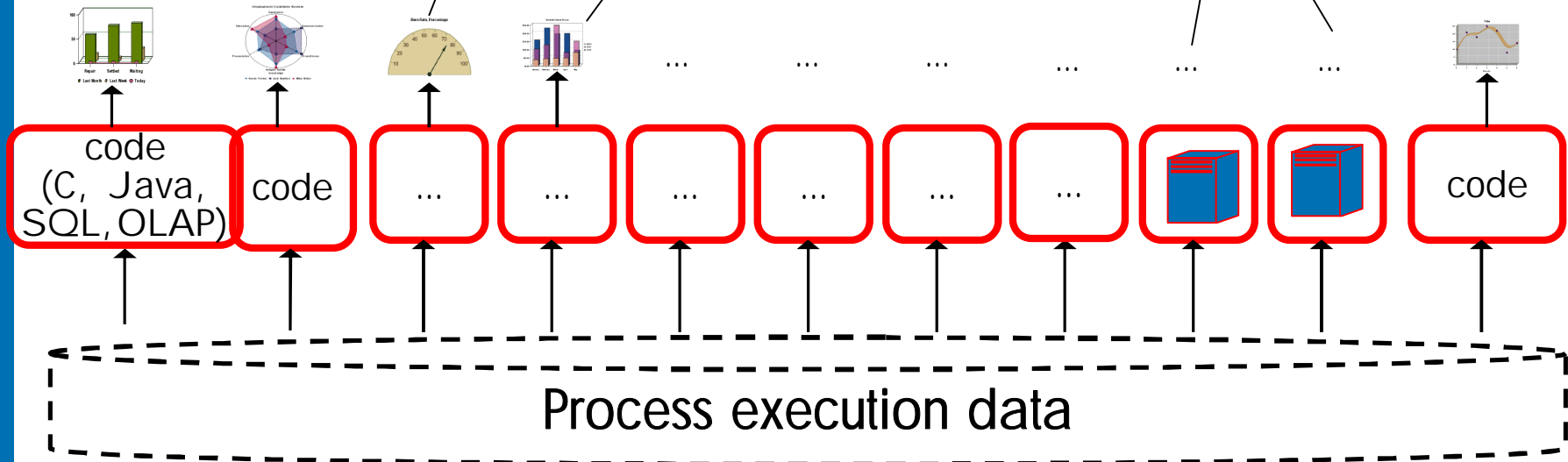average execution quality by process

total 'time to acceptance' SLA violations by customer

predictions and explanations on SLA violations

total execution quality by day

details on total 'time to acceptance' SLA violations

-On time
-Late
-Fail



REPORTS

| Exec quality | SLA violation | outcome | ... | ... | ... | ... |

METRICS

| code | ... | ... | | |

MAPPINGS

CONTEXTS (processes)

27

# Modeling Abstractions

- Domain Model for Processes
  - Entities (e.g., "process", "event", "action", "service")
  - Attributes, relationships, aggregations,..
  - Benefits:
    - Allows definition, computation, analysis, monitoring of "things"
    - Enable easy and quick "verticalization"
- Metric Model
  - Mappings: Functions that compute values from raw data
  - Metrics: Measurable properties of an entity (e.g., transaction value by type) defined by mappings
  - Benefits:
    - Polymorphism (different definitions for different contexts)
    - Minimize number of functions needed
    - Reuse: share definitions across metrics, contexts, data models
- Reporting/Analytic Model
  - Once a metric has been defined, lots of report types are immediately available without requiring coding
    - Domain-specific aggregations
    - Temporal aggregations
    - Analytics on metrics: correlations, predictions, explanations, root cause analyses, etc.
  - Benefits:
    - Rich reporting of generic measures
    - Flexibility: Enabling aggregations of "something" by "something else

# Adding Semantics:
# Behaviors, Contexts, Taxonomies

- Associate semantics with process instances
- Map execution data into categorical metrics
- Examples:
  - Instances lasting more than avg+stdev are slow
  - Orders >$1000 are large orders
- Many behavior templates are predefined
  - T1: "instances in which node N was executed"
- Users define behaviors by instantiating a template (filling forms)
  - T1 (N="notify acceptance") maps to accepted
- New templates can be defined (e.g., in SQL)
- Contexts define the instances to which the template should be applied
  - Processes of group supply chain where suppliers are in group "IT suppliers"
- Define Taxonomies based on behaviors
  - Duration taxonomy: acceptable, fast, slow, unacceptable.

# Sharing mapping functions

Visualization of metric values

| Metric 1 | Metric 2 | ... | ... | ... | ... | ... | ... | ... | ... | Metric n |
|---|---|---|---|---|---|---|---|---|---|---|

Business and IT metrics (e.g., cost, quality, value, performance)

| SQL Query 1 | SQL Query 2 | ... | ... | ... | ... | ... | ... | ... | ... | SQL Query n |
|---|---|---|---|---|---|---|---|---|---|---|

Queries that map execution data into metrics

Process data

| Metric 1 | Metric 2 | ... | ... | ... | ... | ... | ... | Metric k |
|---|---|---|---|---|---|---|---|---|

| mapping 1 | mapping 2 | ... | mapping m<<n |
|---|---|---|---|

Process data

# Development Time: Automatic Query Generation

- Benefits

  - Make things simple, minimize user errors

  - Automatic query generation



```
SELECT EI.EVENT_INST_UUID
FROM EVENT_INSTANCE EI
WHERE (ei.time_completed-ei.time_triggered)< $THRESHOLD$
```

```
SELECT MTE.METRIC_ID, mte.id, EI.EVENT_INST_UUID,
mte.metric_class_id, NULL,NULL, NULL, SYSDATE
FROM active_METERS MTE, CONTEXTS CTX, EVENT_INSTANCE EI,
ENTITIES E, EVENT EV,
EXTENDED_METER_PARS emp
WHERE MTE.MAPPING_ID=5
AND E.EXTENDED_NAME='Event Data'
AND EI.EVENT_ID=EV.ID AND CTX.METER_ID=MTE.ID AND
emp.meter_id=mte.id
AND (ei.time_completed-ei.time_triggered)< emp.PAR1
AND  ((EV.NAME=CTX.EVENT_NAME OR CTX.EVENT_NAME IS NULL)  )
```

# Runtime: Multi-query processing

- Benefits
  - Minimize the number of queries to be executed
    - One mapping can compute several metrics
  - Enable real-time computation and optimizations
    - Smart refresh
    - Take advantage of shared mappings
  - Streaming data, continuous queries

# Correlations between Metrics:
# Root cause & business impact analysis

- Learn relationships between IT data/events and business metrics

| | | | |
|---|---|---|---|
| LOB | Detect no orders placed on any suppliers in the last hour. $$$$$$$$ | Why? | Implies |
| Business Solutions | PO Collaboration solution is not generating orders. | Why? | Implies |
| Applications & Services | I2 business application received no PO requests in the last hour. | Why? | Implies |
| Middleware | EAI bus transmitted no requests to i2.<br><br>EAI / SAP adapter is down. | Why?<br><br>Why? | Implies<br><br>Implies |
| Systems<br>Storage, Devices, Servers | Integration server is down for the last hour. | !!!! | !!!! |

# Outline

- Context
  - Business Process Management
  - Business Process Intelligence
  - Relevance of Information Integration
- Modeling Issues
  - Process Views
  - Metrics Model
- **Information Integration Issues**
  - Generic Data Warehouse Schema
  - Abstraction Mechanisms
  - Generic ETL
  - Information Extraction from Semi-Structured Data
- Summary

# Process Data Warehouse Design: Constellation Schema (almost)

| | |
|---|---|
| Process Defs, process groups | |
| Data Items | Process Behaviors → Behavior defs |
| Time | Process Instances → Resources, Resource groups |
| Node Defs | Task Instances → Services, Service groups |

Must be generic for domains, processes, resources, etc., and yet easily customizable and extensible to new process types, data sources, metrics, report types

Several tricky modeling issues

Also, challenges in how to deal with real-time data, event streams, text, etc.

# Generic Process Warehouse Model

- Challenges for a generic model
  - Multi-level instance data
    - Step level facts, process instance level facts, data-related facts
    - Facts may have to be self-correlated
  - Business data complexities
    - Different from process to process
    - Complex structures
    - Can change at every step during the process
    - à representation hard to generalize
  - Process and step executions go through a lifecycle
    - Step status changes (created, activated, completed, etc --> process events mark progression); number of states can be unlimited (suspend/reactivate cycles)
    - Different systems supporting the execution may have different lifecycle phases

# Main elements of the generic warehouse model

- Single granularity for steps (rather than at the level of status changes)
- Single fact table for any step of any process
  - Enables analyses across processes
  - Includes aggregation of most common step event measures
- Correlation with previous step data handled via additional columns
- Separate business data tables for each process type
- Blind links to handle step/process correlation with business data

# Process warehouse model

# Mapping events to abstract processes

- Two facets to provide abstracted process representations

  - A way to model the abstraction

    - Describe the high level process
    - Describe how its progression maps to underlying IT events

  - ETL mechanism to load warehouse with abstracted process execution data

# Map Events to Abstract Process Progression

Implemented process

- Typically, events signal status changes in steps of the implemented process
- Have to specify or learn abstract process progression
- Mappings between monitored events and start/completion of abstract process steps, data relevant to the abstract process, …

Audit Msg
Invoice_ID=123
Amount=$100
response=OK

Event

Process data changes map to progression information



Invoice received → Image import → Work object creation / Indexing → Index validation → PO correction → Automatic audit / Manual audit → Done / Correction → Reject invoice

Abstract invoice payment process

# Two-phased mapping

Audit Msg
Invoice_ID=123
Amount=$100
response=OK

IT Event captured by a probe

Invoice (abstract data)

Map event to business data changes

map business data changes to process progression information

Invoice received

Image import

Work object creation

Indexing

Index validation

PO correction

Automatic audit

Manual audit

Correction

Done

Reject invoice

# Why the indirection?

- Many different events may cause the same change to a business data item

- Same business data can be used to support and mark progression of instances of different process types

- In practice, for abstract processes the progression often depends on business data changes

- Benefits

  - Reduces specification & maintenance effort

  - Specs are more robust to changes in the info sources (event specs updated but no need for business data or progression info)

# Extraction & abstraction of process data

Process progression and business data ← Process mapping ← **Process to data mapping defs**

**Staging Area**

Business data changes ← Business data mapping ← **Data to event mapping defs**

Event data ← Data extraction

Event log ← probe ← Web server

Event log ← probe ← Application server

ERP log ← probe ← ERP system

Message broker log ← probe ← Message broker

Event log ← probe ← Sys mgmt tool

# Loading process data

- Modeling specs used by the ETL to map across levels of abstraction

  - IT events captured with probes and logged with timestamps

  - ETL reads event tables in logs and orders them by time

  - Events are mapped to business data changes

  - Business data changes are 'replayed' in order and relevant changes are detected for computing process progression

  - Process progression creates records for the step execution data which are loaded into the warehouse

# ETL generation

- Automates staging area creation & maintenance

- Automates generation of executable transformation scripts

  - Indirection of mappings from IT events to process progression à Two-phased transformation

    - Phase 1: IT events mapped to business data changes
    - Phase 2: business data changes mapped to process progression

# Staging area

- Three types of tables
  - Landing tables
    - Buffering of extracted IT events data
    - Checks for errors in the extraction
    - Refreshes at every cycle
  - Image tables
    - Keep an image of the IT events records extracted since the first extraction
    - Input to first transformation phase
  - Comparisons between landing & image tables
    - To detect duplicates
    - Determine manipulation operation (I, d, u)
  - Intermediate tables
    - Output of first transformation phase
    - Business data changes
    - Input to second transformation phase

# Intermediate tables

- Alternative design: 2 separate ETL processes but …

  - Inefficient

    - Extraction and staging of business data changes
    - Additional tables to keep all business data changes to mark process progression

      - DW only stores the last version of a business data instance

# ETL Transformation phases

Extract

Transformation
Phase 1

Transformation
Phase 2

Load

Staging Area

Log 1

E

Landing
tables

Image
tables

Inter-
mediate
tables

Process
DW

Log n

Mapping Generation

IT event-
Biz data
mappings

Biz data-
process
exec data
mappings

BPI Repository

# Executable mapping generation

- How to execute the transformations?
  - Agnostic to underlying tool
  - Modeling: declarative mappings
  - Mapping Generator derives prescriptive mappings
    - Two phases
      - Prescriptive logical mappings
        - Canonical language to express executable semantics (pseudo-SQL)
      - Prescriptive executable mappings
        - Specific translators (or manually)
      - Orthogonal to the two transformation phases

# Mapping Generator

- Core: mapping templates
  - Parameterized logical scripts in canonical language
    - Capture executable semantics
      - Factor out commonalities of mapping between the layers of abstraction
      - Exploits DW semantics
      - Captures other correspondences not specified by the declarative mapping (e.g., duration)
    - Parameters: event-, business entity-, process step-related
    - Templates instantiated by declarative mappings
    - Different template types (e.g., bizEntity_to_endStep)
    - Not executable
    - Canonical language translator

# Mapping generation phases

# Outline

- Context
  - Business Process Management
  - Business Process Intelligence
  - Relevance of Information Integration
- Process Modeling Issues
  - Process Views
  - Metrics Model
- Information Integration Issues
  - Generic Data Warehouse Schema
  - Abstraction Mechanisms
  - Generic ETL
  - **Information Extraction from Semi-Structured Data**
- Summary

# Information Extraction from Less Structured Data

- Information relevant to business processes may be contained in less structured data

- Many enterprises have a large corpus of contracts, customer service logs, reviews, etc.

- For example, the enterprise must be able to respond to events that might affect existing contractual relationships

- Critical information remains buried in text, e.g., key parameters of Service Level Agreements

- Need to incorporate this information into Business Process Intelligence solution

# Our approach

- Automatically identify "facts" in text documents
- Based on the use of
  - two object models
  - information extraction techniques
- Identified facts can be
  - automatically tagged with XML
  - extracted into the business process data warehouse as additional business data
- Extracted information becomes readily available for BPI: metrics definitions, queries, reports, analyses

# Document templates

- Common practice: set of free text templates for different contract types
  - For each type à one or more templates
  - E.g.,
    - Long term agreement (LTA)
    - Corporate purchase agreement (CPA)
- Contract templates organized in clauses
  - Each clause: specific kind of factual information
  - E.g.,
    - Term clause

# CPA term clause template

This CPA will be a [TERM] Agreement for the period [START DATE] to [EXPIRATION DATE] inclusive. Both parties agree to meet prior to [MM/DD/YY] to consider an extension of [##] year(s). In like manner, both parties shall meet prior to [MONTH/DAY OF EXPIRATION DATE] of each year to consider future extensions.

# Template instantiation

- Relevant fact types in templates: attributes
- Clause template instantiation
  - Value assignment to attributes (facts)
- Contract
  - Combination of clause template instantiations
  - Some variations
    - In the text
    - In the order of clauses

# Instance of CPA term clause

This CPA will be a one year Agreement for the period 05/01/03 to 05/01/04 inclusive. Both parties agree to meet prior to 04/01/04 to consider an extension of one year. In like manner, both parties shall meet prior to 05/01 of each year to consider future extensions.

# Document annotation

- Learning rules to automatically identify facts requires a training set
  - Subset of contract collection annotated with tags for relevant facts
- Two kinds of tags
  - Attribute tags
  - Structural tags

# Attribute tags in term clause instance

This CPA will be a *<TERM>* one year *</TERM>* Agreement for the period *<START_DATE>* 05/01/03 *</START_ DATE>* to *<EXPIRATION_DATE>* 05/01/04 *</EXPIRATION_DATE>* inclusive. Both parties agree to meet prior to *<IMMEDIATE_EXTENSION_MEET_DATE>* 04/01/04 *</IMMEDIATE_EXTENSION_ MEET_DATE>* to consider an extension of *<EXTENSION_PERIOD>* one *</EXTENSION_PERIOD>* year. In like manner, both parties shall meet prior to *<FUTURE_EXTENSION_MEET_DATE>* 05/01 *</FUTURE_EXTENSION_MEET_DATE >* of each year to consider future extensions.

# Structural tags for term clause

*<TERM_CLAUSE>* This CPA will be a <TERM> one year </TERM> Agreement for the period <START_DATE> 05/01/03 </START_ DATE> to <EXPIRATION_DATE> 05/01/04 </EXPIRATION_DATE> inclusive. Both parties agree to meet prior to <IMMEDIATE_EXTENSION_MEET_DATE> 04/01/04 </IMMEDIATE_EXTENSION_ MEET_DATE> to consider an extension of <EXTENSION_PERIOD> one </EXTENSION_PERIOD> year. In like manner, both parties shall meet prior to <FUTURE_EXTENSION_MEET_DATE> 05/01 </FUTURE_EXTENSION_MEET_DATE > of each year to consider future extensions *</TERM_CLAUSE>*

# Object models

- Objective: guide annotation task
- XML-based
- Two kinds
  - Document object model (DOM)
    - Specifies structural components
      - Sections & clauses
      - Order can vary
  - Facts object model (FOM)
    - Specified relevant attributes
      - E.g., contract expiration date

# Document Object Model

```
<DOM>
        <id>  7002 </id>
        <contract type>
                LTA
        </contract type>

                …
        <section>
                <name> Shipment and Delivery </name>
                <clause> prospective failure </clause>
                <clause> untimely shipment </clause>
        </section>
        <section>
                <name> Term </name>
                <clause> term </clause>
        </section>

                …
</DOM>
```

# Facts Object Model

```
<FOM>

        <id> 235 </id>
        <contract type>
                LTA
        </contract type>
        <attribute>
                <name> expiration_date </name>
                <type>  date </type>
                <nature> mandatory </nature>
        </attribute>
        <attribute>
                <name> untimely_transportation_
                        means</name>
                <type>  transportation </type>
                <nature> mandatory </nature>
        </attribute>
            …
    </FOM>
```

# Semantic types

- Enumeration

    &lt;type&gt;

        &lt;name&gt; transportation &lt;/name&gt;

        &lt;kind&gt; enumeration &lt;/kind&gt;

        &lt;values&gt; airplane, ship, truck, trailer &lt;/values&gt;

    &lt;/type&gt;

- Format

    &lt;type&gt;

        &lt;name&gt; date &lt;/name&gt;

        &lt;kind&gt; format &lt;/kind&gt;

        &lt;values&gt; mm/dd/yy, month dd year, mm-dd-yyyy
                                    &lt;/values&gt;

    &lt;/type&gt;

# Regularities

- Each attribute (fact type) can have one or more associated regularities

- Structural regularity
  - Regularities in the structural component (location) of an attribute
  - E.g., untimely_transportation_mean à  Shipment and Delivery section

- Phrasal regularity
  - Regularities in the surrounding words
  - E.g., for the start_date attribute of a term clause
    - for the period 01/01/2004 to
    - starting from 01/01/2004 "until

- Grammatical regularity
  - Regularities in the parts of speech (e.g., noun, verb, adjective, etc) of surrounding words, and/or in the syntactic relations between them (subject, etc)
  - Take advantage of clausal structure provided by a syntactic analyzer and PoS tagger

# Extracting facts

- Text mining
  - Information extraction techniques
    - Learn patterns for regularities
    - Train on annotated set
    - Generate rules
      - Antecedent: pattern for combination of regularities
      - Consequent: attribute name for corresponding fact
    - One or more rules for each attribute
      - Different kinds of contract types
      - Different contract templates for a same contract type
      - Text variations
    - Rules database

# Rule example

- Rule for attribute expiration_date from an LTA term clause instance

```
<Rule>
    <id> 153 </id>
    <FOM_object> 235 </FOM_object>    //id for LTA FOM object
    <antecedent>
            <structural_component>
                    <section> TERM </section>
                    <clause> TERM </clause>
            </structural_component>
            <surrounding_component>
                            'period' date 'to' (date)
            </surrounding_component>
    </antecedent>
    <consequent>
            <attribute> expiration_date </attribute>
    </consequent>
</Rule>
```

# Generating rules

- Information extraction technique
  - Each instance-attribute tag pair becomes a "seed" to grow a new rule that covers the seed
  - Top-down algorithm to induce rules
    - First finds the most general rule that covers the seed
      - Anchors the boundaries of the fact
    - Then, extends the rule by adding terms one at a time
      - Metric to select a new term: expected error of the rule
- The technique is made more efficient by the use of structural tags
  - At rule generation time the structural context narrows the validation space
  - At rule application time the structural context narrows the search space

# Summary

- The intelligent enterprise monitors and optimizes its business processes and interactions with business partners.

- Better business process management is "essential" (and independent of automation)

- Today, this is very difficult to do, although tools are appearing to address pieces of the problem.

- Our approach: a Business Process Intelligence solution (Business Cockpit) that combines process modeling, metrics definition, generic DW schema and ETL generation, and analytics

- But, many research challenges remain.

# References

- Fabio Casati, Malu Castellanos, Umeshwar Dayal, Norman Salazar,"A Generic Solution to Warehousing Business Process Data." VLDB 2007.
- Fabio Casati, Malu Castellanos, Norman Salazar, Umeshwar Dayal, "Abstract Process Data Warehousing." ICDE 2007.
- Fabio Casati, Malu Castellanos, Umesh Dayal, Ming-Chien Shan, "A Metric Definition, Computation, and Reporting Model for Business Operation Analysis." EDBT 2006.
- Ming C. Hao, Daniel A. Keim, Umeshwar Dayal, Jörn Schneidewind, "Business Process Impact Visualization and Anomaly Detection" Information Visualization Journal 2006.
- Malu Castellanos, Fabio Casati, Mehmet Sayal, Umeshwar Dayal, "Challenges in Business Process Analysis and Optimization." Proc. TES Workshop, Springer-Verlag, 2005.
- Malu Castellanos, Fabio Casati, Umesh Dayal, Ming-Chien Shan, iBOM: A Platform for Business Operation Management." ICDE 2005.
- Fabio Casati, Malu Castellanos, Umesh Dayal, Ming-Chien Shan, "Probabilistic, Context-Sensitive, and Goal-Oriented Service Selection." ICOSOC 2005.
- Malu Castellanos, Norman Salazar, Fabio Casati, Umesh Dayal, Ming-Chien Shan, "Predictive Business Operations Management." DNIS 2005.
- Mehmet Sayal, Ming-Chien Shan, "Analysis of Numeric Data Streams at Different Granularities." IEEE International Conference on Granular Computing, July 2005.
- Malu Castellanos, Fabio Casati, Umeshwar Dayal, Ming-Chien Shan, "A Comprehensive and Automated Approach to Intelligent Business Process Execution Analysis." Distributed and Parallel Databases 16(3): 239-273, 2004
- Malu Castellanos, Norman Salazar, Fabio Casati, Ming-Chien Shan, Umesh Dayal, "Automatic Metric Forecasting for Management Software." OVUA Workshop 2004.
- Daniela Grigori, Fabio Casati, Malu Castellanos, Umesh Dayal, Ming-Chien Shan, Mehmet Sayal. "Business Process Intelligence." Computers in Industry 53 (3), April 2004.
- Ming C. Hao, Daniel A. Keim, Umeshwar Dayal: VisBiz, "A Simplified Visualization of Business Operation." IEEE Visualization 2004
- Malú Castellanos, Fabio Casati, Umeshwar Dayal, Ming-Chien Shan, Intelligent Management of SLAs for Composite Web Services. DNIS 2003.
- Fabio Casati, " Eric Shan, Umeshwar Dayal, Ming-Chien Shan:, "Business-oriented management of Web services." Commun. ACM 46(10)
- Fabio Casati, Umeshwar Dayal, Ming-Chien Shan, " Business Operation Intelligence." DNIS 2002.
- Mehmet Sayal, Fabio Casati, Umeshwar Dayal, Ming-Chien Shan, "Business Process Cockpit." VLDB 2002.
- Angela Bonifati, Fabio Casati, Umesh Dayal, and Ming-Chien Shan, "Warehousing Workflow Data: Challenges and Opportunities." VLDB 2001.

# Thanks!