# Real-Life Data Mart Design Challenges

**Leslie Tierstein**

**newScale®**
*Service Catalog Leader*

# Overview

- Data mart design challenges addressed in this project:
  - Dynamically defined dimensions
  - Potentially 100's of dimensions
  - Loading XML data into the data mart
  - An incremental refresh that includes updates
  - Supporting multiple databases with a standard toolset and code base
  - Designing the physical database to map to the BI tools' business view

- Data mart design challenges not addressed in this project:

  - Performance tuning for very large databases (VLDB's)

  - Data cleansing and data integrity

  - Integrating data from heterogeneous databases or multiple applications

- Service Catalog Management
- An IT department defines a catalog of all the services it provides to users
  - End User Computing: Add Visio
  - Advanced Computing: Configure a server
- Users request a service
  - System tracks the service request, all tasks required to complete it

# The OLTP System (1)

● Requestor View – Upgrade Memory

   – User fills out and submits an order (requisition)

# The OLTP System (2)

- IT View – Upgrade Memory
  - User request has criteria for completion
  - Request is routed to appropriate people for authorization and task delivery

**Services**

| Name | Service Level Description | Standard Duration[+] | Quantity | Unit Cost | Subtotal |
|------|--------------------------|---------------------|----------|-----------|----------|
| Computer Memory - Upgrade | | 3 business days | 1 | 300.00 | 300.00 |
| | | | | **Total Cost:** | **300.00** |

[+] Standard Duration applies to delivery after any required authorizations have been completed.

**Delivery Process**

| Process Milestone | Due Date | Completed On | Status |
|-------------------|----------|--------------|--------|
| Service Group Authorization | 12/28/2006 4:00 | | In Progress |
| Delivery project for Computer Memory - Upgrade | 01/02/2007 10:00 | | Pending |

● Service Designers' View

- – Standard information (duration, cost, delivery plan, authorizations) defined via UI
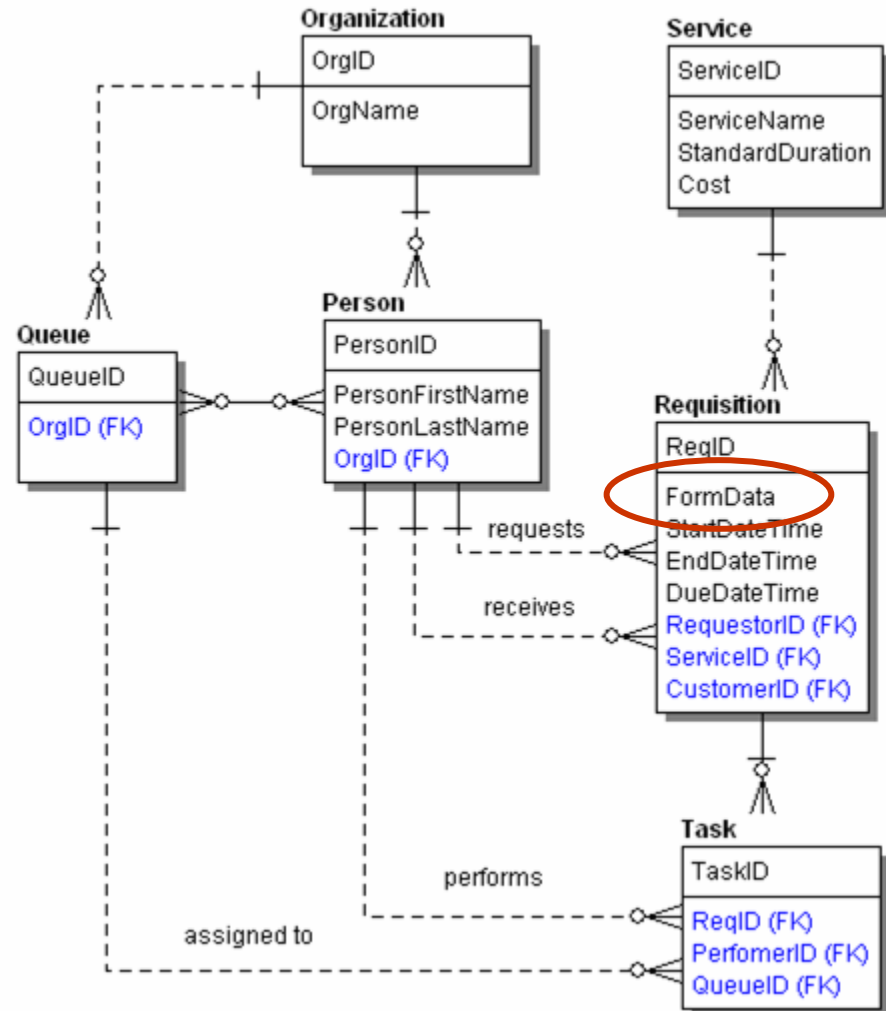- – Service-specific data entry requirements defined via "dictionaries"

- Vastly simplified ERD of transactional data

- Metadata on service definitions is also maintained

- Where is form (dictionary-based) data?

# Data Mart Requirements

- Data mart must include form data

- Refresh needs to run in nightly maintenance window

  – Minimize impact on global accounts

- Data mart must include data on both open (in progress) and completed requisitions and tasks

  – Refresh must perform "updates"

- Extract XML data from the requisition BLOB
  - Cognos8 includes Composite, which can extract data from XML
    - Performance is problematic
  - Adding a custom plug-in in the ETL tool provided insufficient power/flexibility
  - Custom Java program was needed
- Would the custom program perform acceptably?

# ETL Performance Studies

- 100,000 requisitions per month (an order of magnitude higher than current usage)

| Data Points | # | Units |
|---|---|---|
| Average # of requisitions started per day | 5000 | Requisitions |
| Average # of entries per requisition | 1.5 | Entries per Requisition |
| Days from requisition start until it's closed | 10 | Days |
| Average # of tasks in a service | 6 | Tasks per Service |
| Average # of actions taken per task | 1.6 | Actions per Task |
| Average # of dictionary data elements | 10 | Data Elements per Dictionary |
| Average # of dictionaries per service | 4 | Dictionaries per Service |
| % requisition entries with reportable dictionaries | 70% | Requisition Entries |
| Time to extract one data element to data mart | 8 | Milliseconds |

| Analysis | | Units |
|---|---|---|
| Requisitions open at any one time | 50,000 | Requisitions |
| Requisitions entries open at any one time | 75,000 | Requisition Entries |
| Tasks open at any one time | 45,000 | Tasks |
| Actions recorded per day | 72,000 | Actions |
| Requisitions entries changed in any one day | 72,000 | Requisition Entries |
| Requisition entries with reportable dictionaries changed per day | 50,400 | Requisition Entries |
| Requisition entry dictionaries changed per day | 201,600 | Requisition Entry Dictionaries |
| Requisition entry data elements changed per day | 2,016,000 | Requisition Entry Data Elements |
| **Time to extract changed data elements** | **16,128** | **Seconds** |
| | **269** | **Minutes** |
| | **4.48** | **Hours** |

- Averages are good, but know your outlyers
  - If only to be able to warn users about potential issues
  - Hotel/Hospitality Company:
    - 30% of services have 100-160 fields, with a maximum of 250. About 60% have 70-100 fields
    - There are dictionaries with about 40 fields, with a maximum of 43.

- Inserts are easy, updates … NOT

- Data Mart contents need to be updated
  - Tasks and requisitions are added to the data mart as soon as the requisition is submitted
  - Status of requisitions may change, as well as completion date and form data
  - Task status and info is updated as the task is performed (completed)
- Logical "update" <> DML UPDATE

- Argh! Updating data mart contents is expensive

  - Updates require indexes

  - Inserts perform much better without indexes

  - Deletes (from large tables) are a disaster

- Avoid updates by placing data to be updated in a partition that is truncated

  - Refreshed data is written to the same partition or to a permanent partition when complete

- Still an issue

- Partitioning not implemented because:

  - Small table size and growth rates

  - Lack of support in ETL tool

  - SQLServer (partition views only!)

  - What is partition key?

- Using insert/update model

- Time will tell if this is robust enough

- Table, View or Metadata?
  - Design the dimensional model in conjunction with the business view of the model
    - Date dimension(s)
    - People dimension(s)

- Date required as a dimensional key
  - Due Date – for task and requisition
  - Closed Date – for task and requisition
  - Start Date – for task and requisition

- Complete date and time may also be needed
  - SLAs and OLAs

**Requisition**

| ReqID |
| --- |
| FormData |
| StartDateTime |
| EndDateTime |
| DueDateTime |
| RequestorID (FK) |
| ServiceID (FK) |
| CustomerID (FK) |

**Task**

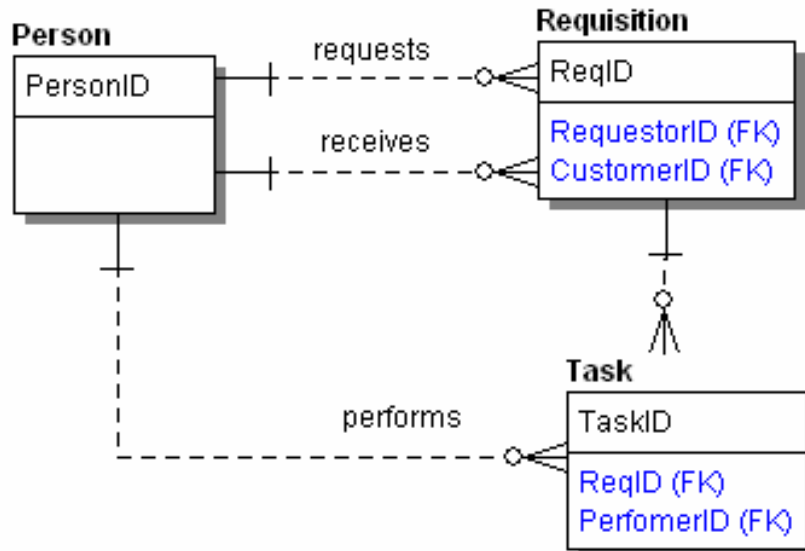| TaskID |
| --- |
| StartDateTime |
| EndDateTime |
| DueDateTime |
| ReqID (FK) |
| PerfomerID (FK) |
| QueueID (FK) |

- How many database objects for the date dimensions?
  - 3 tables, 3 subject areas
  - 1 table, 3 views, 3 subject areas
  - 1 table, 3 subject areas
- Pros and cons of each approach
  - ETL, debugging, user interface in ReportNet
  - Where are relationships defined?
    - In the BI tool
    - In the database

- Date dimensions are densely populated
  - Likely to have an attribute for every weekday date value

- Date dimensions can be statically populated
  - Rather than via fact/dimension build

- Solution: 1 table, 3 views, 3 subject areas
  - In Oracle, constraints can be defined in the view and created in the BI tool

- Required Dimensions
  - Customer – for requisition
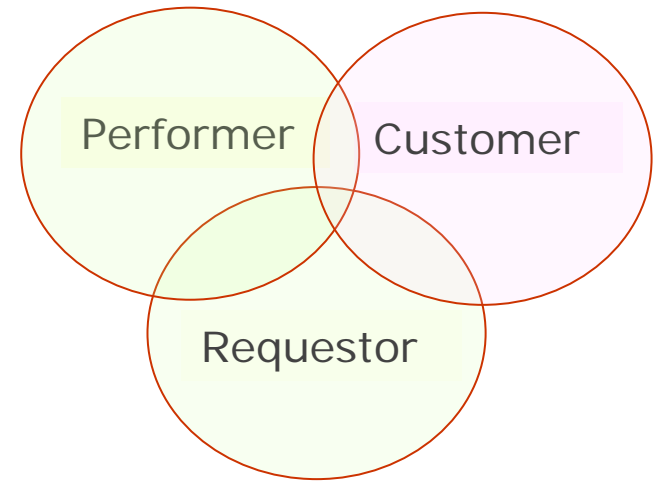  - Requestor – for requisition
  - Performer – for task

# People Dimensions (2)

- **How many database objects?**
  - Do people fulfilling different roles intersect? And how much?
    - Customers and Requestors may overlap 100%
    - Performers are a much smaller group
  - Implications for prompts and filters
    - Prompted filter should only show people who have filled the corresponding role
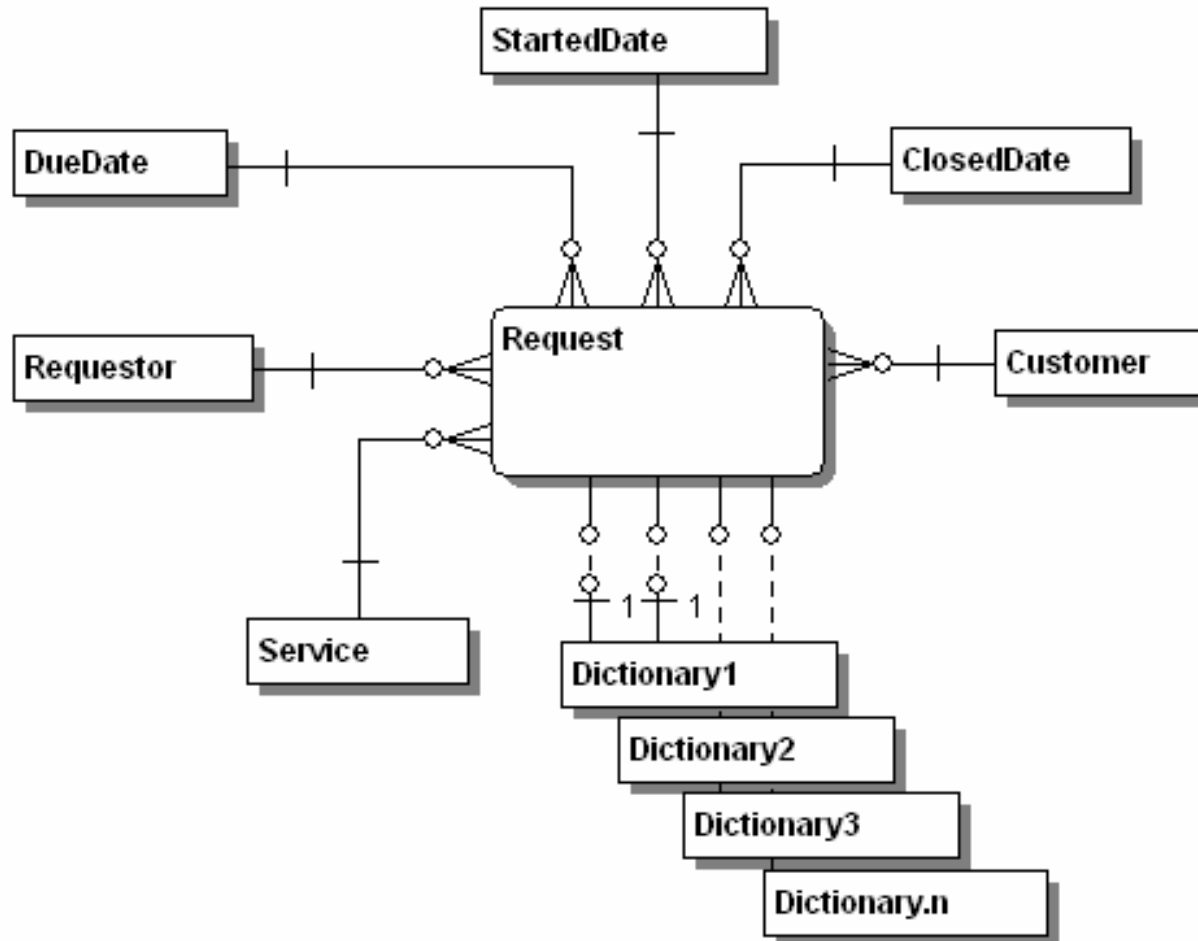
Performer

Customer

Requestor

# Dimensional Modeling

- Dynamic data model to accommodate form data
  - Other companies have done this, but their scenarios are not as complex
    - Oracle Noetix Views and Noetix data warehouse – Flex fields

# Dimensional Model

- How to model the dimensions which correspond to dynamically defined dictionaries
  - Each attribute in a dictionary is a dimension – rejected
    - W-A-A-Y too many dimensions
  - "Junk", "Degenerate" or "Demographic" dimensions
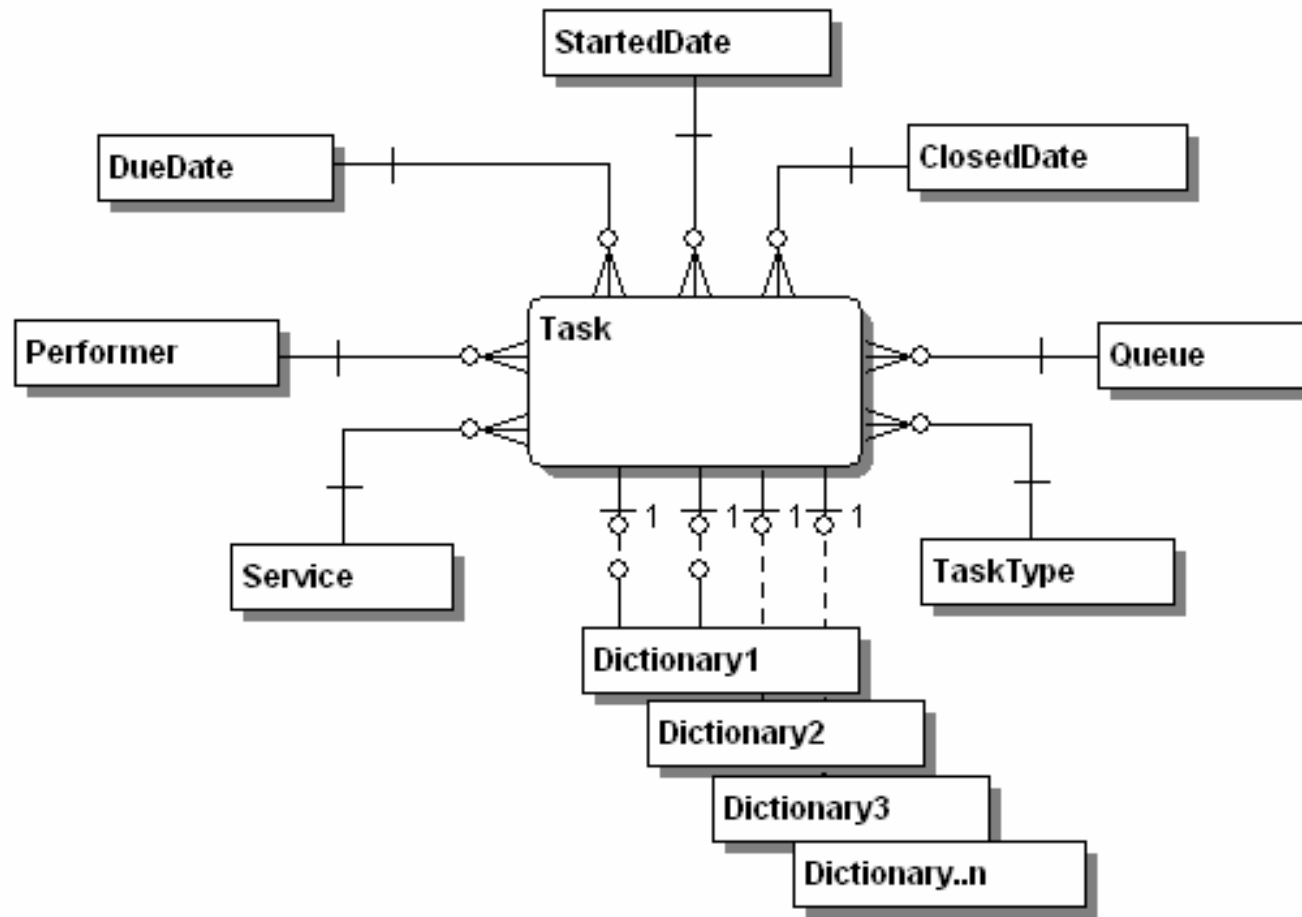  - "Reference" dimensions

- Populate the dictionary dimension using a unique combination of values for each dictionary across all requisitions

- The fact table has a column for a relationship with each dimension – many will be "N/A"

- A complex fact/dimensional build would be required

- See Kimball, Design Tips 46, 48

- Each dimension has a 1:1 relationship with the fact table

- Best to not have basic analytical metrics and frequent query topics in the reference dimension

- Greatly reduces the size of the fact table

- See Kimball, Design Tip 86 (http://www.kimballuniversity.com/html/designtips.html)
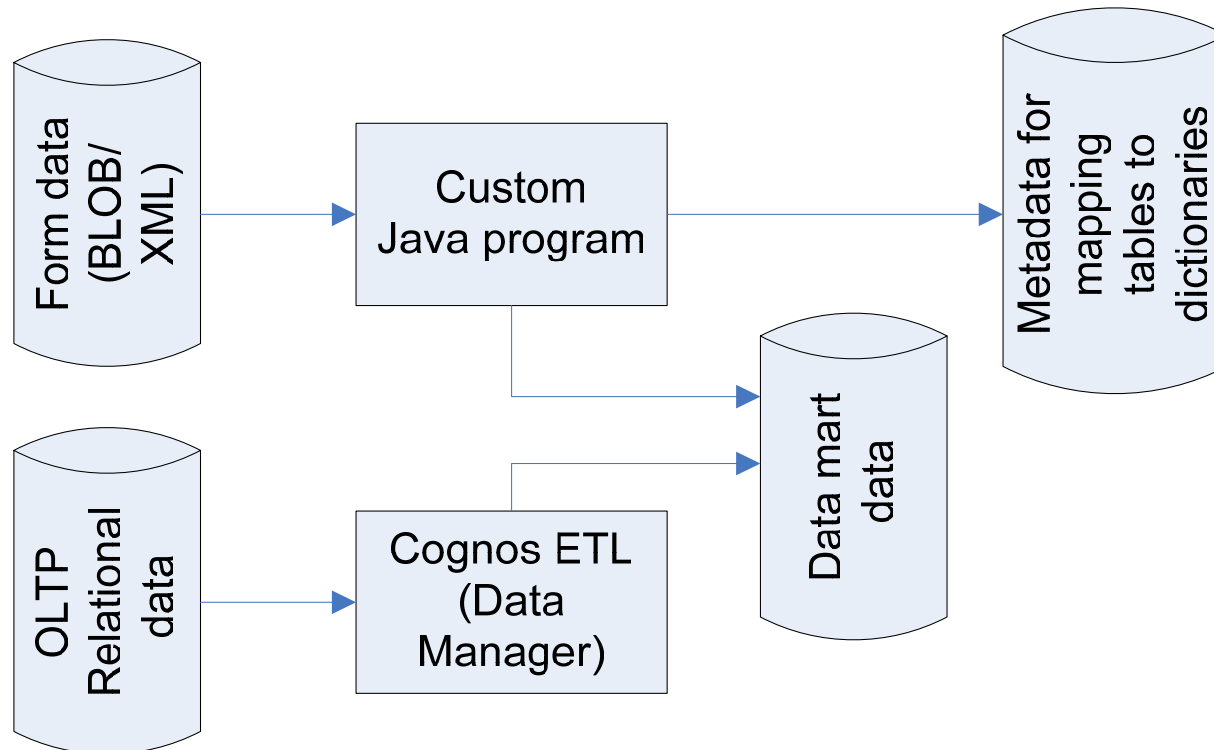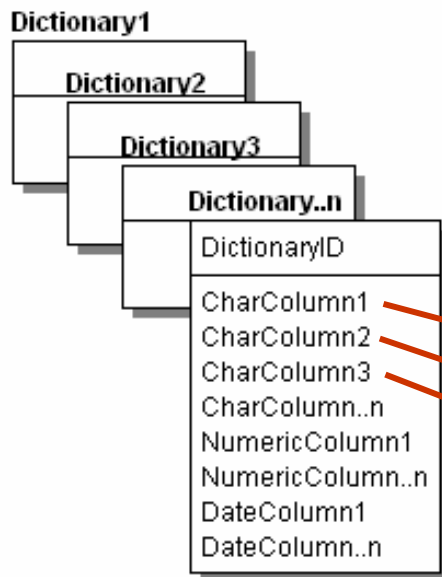
# Task Star Schema

# Tool Suite
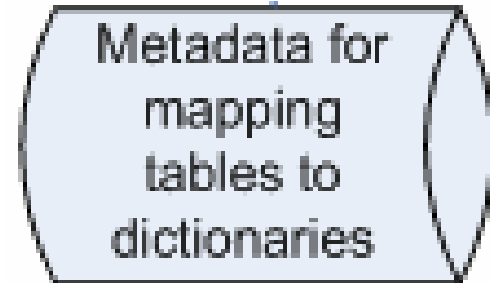
| Tool Type | Cognos Product |
|---|---|
| ETL tool (maps OLTP data to data mart) | Data Manager |
| BI Framework (maps database to end-user/business view) | Framework Manager |
| BI Power User Tool (allows end users to create reports/queries) | ReportStudio QueryStudio |
| BI End User Tool (provides portal for running reports/queries) | ReportNet |

# Data Mart Refresh

- Custom Java program is used in conjunction with Data Manager ETL

# Data Mart Refresh …

● … and metadata creation

Metadata for mapping tables to dictionaries

– Dictionaries and their attributes must be mapped to tables/columns in the data mart

– This mapping is used in the BI business view

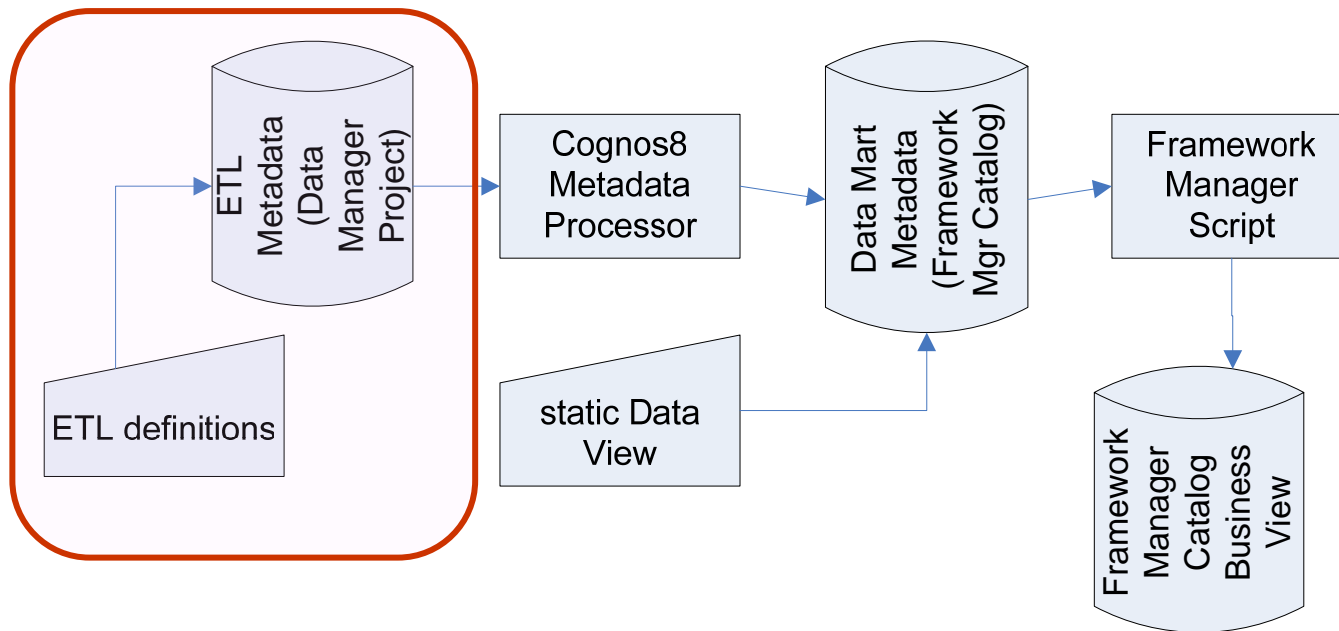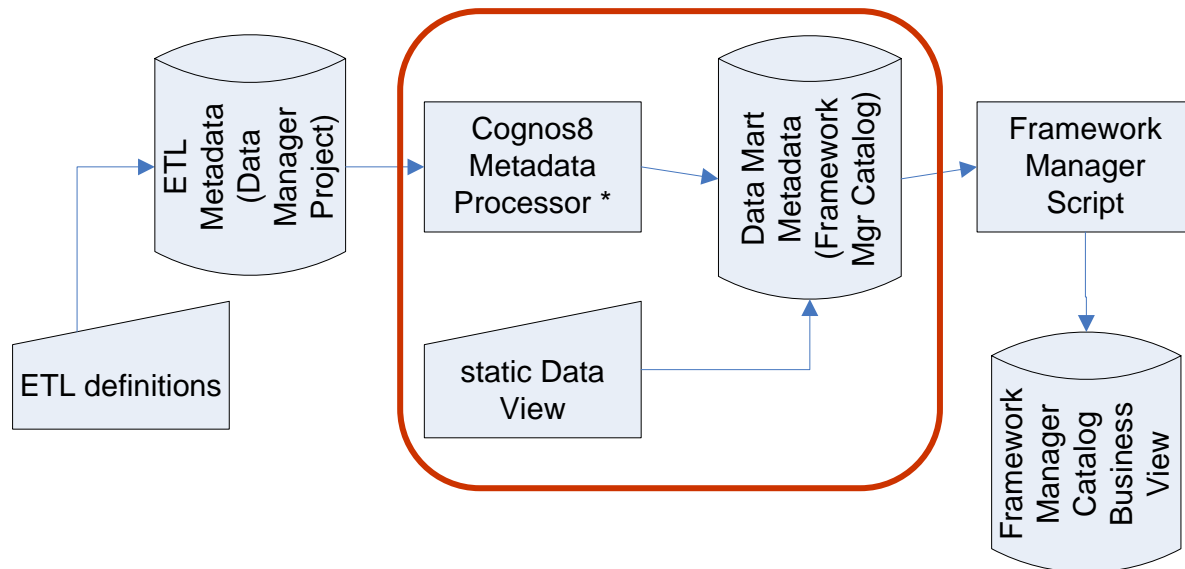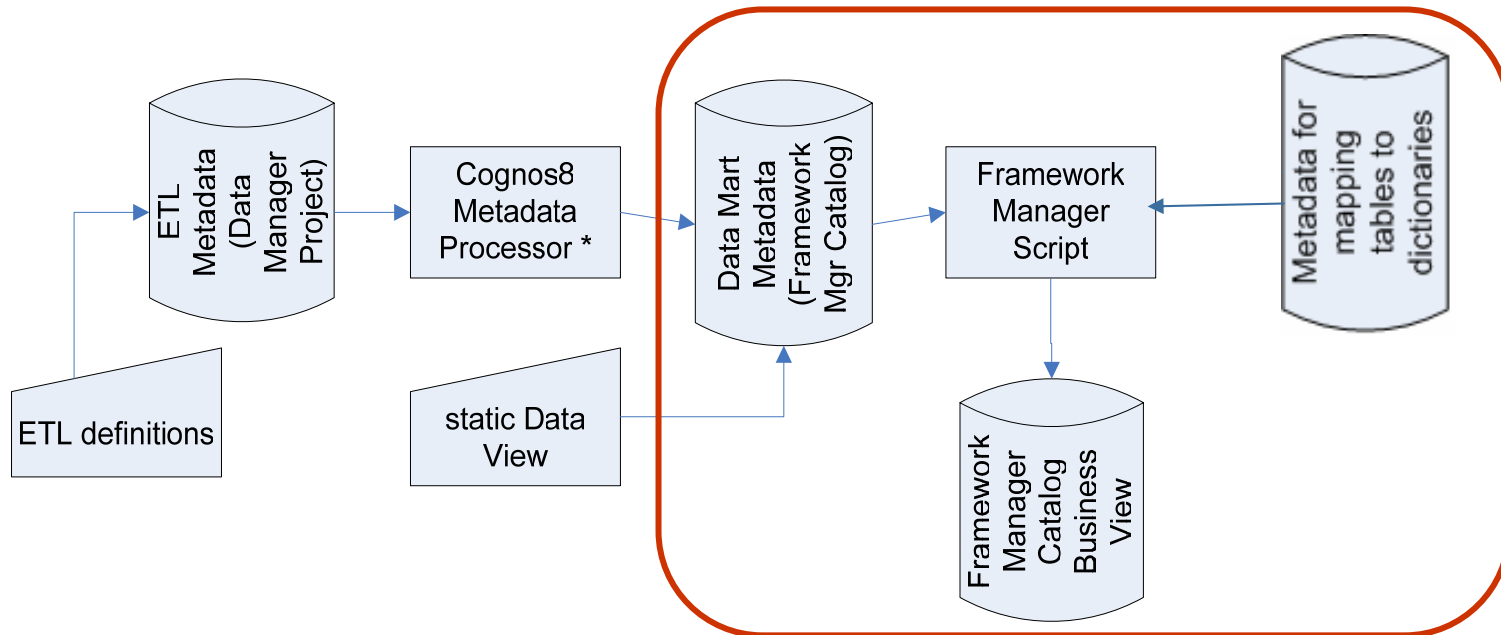- Specify ETL for facts and traditional dimensions

  – Uses Cognos8 ETL tool, Data Manager

- Transform the Data Manager project into a Framework Manager catalog
  - (Sigh) two tools, two different sets of metadata
  - Data Manager's ETL destination tables become the basis of the BI tool's metadata

- For dictionaries ("junk" dimensions) use the Framework Manager API to create a business view
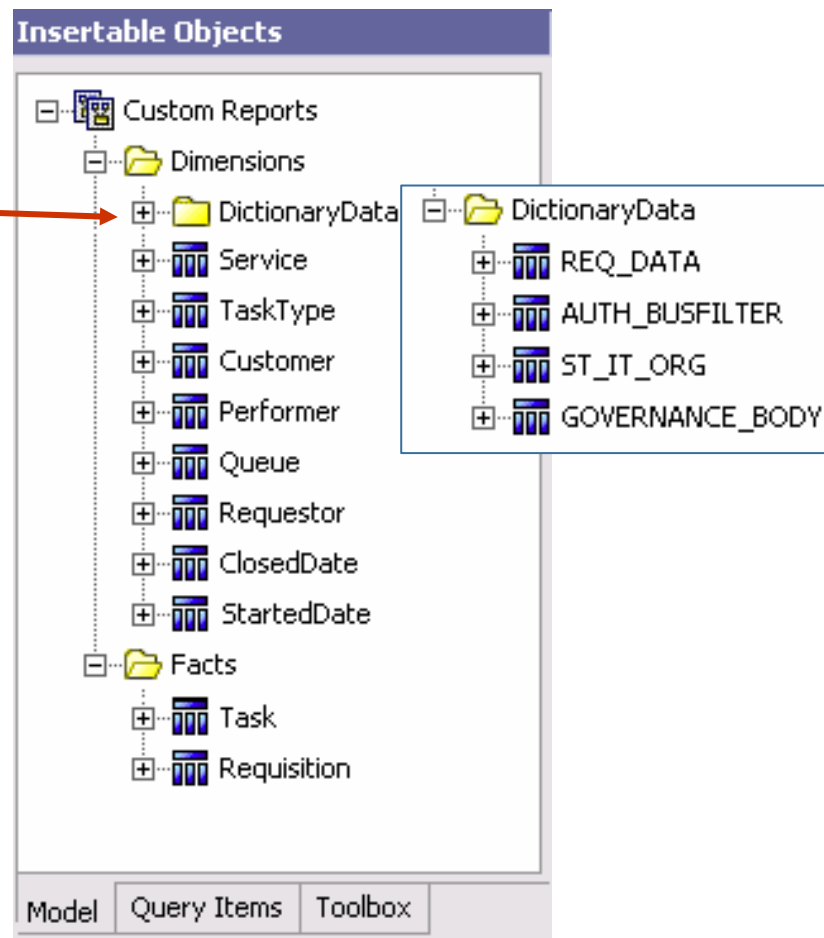  - Appropriate names for dictionary dimensions and data elements

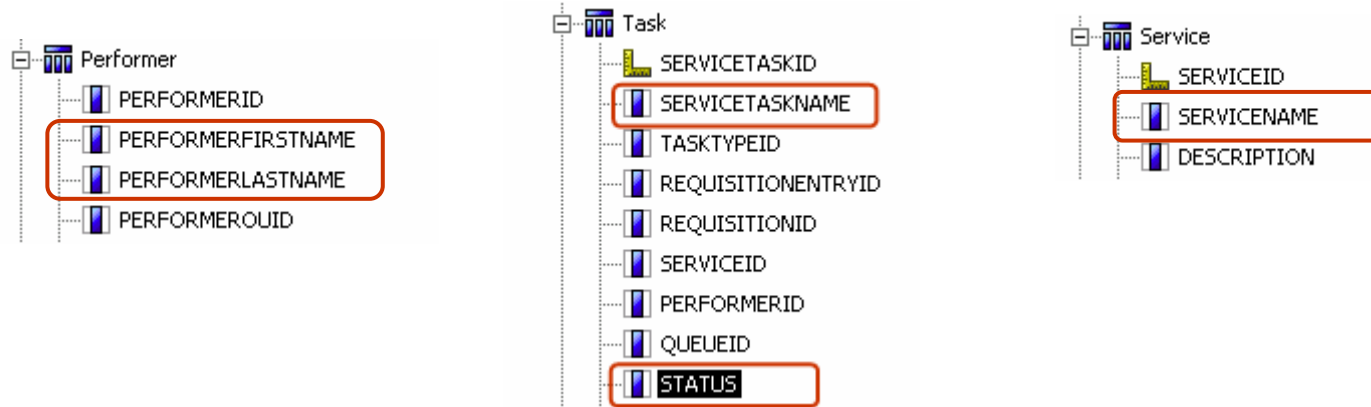# Dimensions

- Dimensions
  - Dynamic dimensions
    - Based on dictionaries designated as reportable
    - Different for each site
  - Universal/static dimensions
    - Found at all sites

● Create a report listing tasks performed, the performer, the service in which they were performed, and the current status

# Conclusions (1)

- XML extraction
  - Java programming was much more efficient than ETL tools available to us

- Dimensional modeling
  - Junk dimensions, reference dimensions, degenerate dimensions – who knew?
  - Subscribe to Kimball
  - Hire short-term consultants

# Conclusions (2)

- ETL Performance
  - Satisfactory
  - Need sophisticated scheduler to track all jobs
- Ad-hoc Report Performance
  - Insufficient testing on the largest user data volumes (IMHO)
  - The software went live on January 31, 2007
  - Ask me in a few months

- **Performance issues**
  - Reporting behavior is acceptable except for building filter selection lists on large fact table
  - Rework was required for customer load at client with largest customer base

- **Database issues**
  - Incompatibilities between Oracle SQL, SQLServer SQL, DB2 SQL, and Cognos SQL
  - Implement in database vs implement in tools

- Functional/design issues
  - Should have gone with two people-based tables
    - Requestor/Initiator and Performer
  - Treatment of dimensions with all blank attributes needs more attention
  - Better use of framework to incorporate display definitions — better metadata required

- Leslie Tierstein is a Senior Technical Architect at newScale, Inc, in Foster City California.

- She can be reached at ltierstein@earthlink.net or leslie.tierstein@newscale.com

- This paper is available on line at: http://home.earthlink.net/~ltierstein