The MIT 2008 Information Quality Industry Symposium

# Using Conceptual Data Modeling to ensure high Information and Data Quality

Pete Stiglich
Senior Consultant
PStiglich@ewsolutions.com

www.EWSolutions.com

Strategic Partner & Systems Integrator

© 2008 Enterprise Warehousing Solutions, Inc. (EWSolutions) – 1

Intelligent Business Intelligence℠

---

# EWSolutions' Background

**EWSolutions** is a Chicago-headquartered strategic partner and full life-cycle systems integrator providing both award winning strategic consulting and full-service implementation services. This combination affords our clients a full range of services for any size enterprise information management, managed meta data environment, and/or data warehouse/business intelligence initiative. Our notable client projects have been featured in the Chicago Tribune, Federal Computer Weekly, Crain's Chicago Business, and won the 2004 Intelligent Enterprise's RealWare award, 2007 Excellence in Information Integrity Award nomination and DM Review's 2005 World Class Solutions award.

**Information Integrity Coalition**

*2007 Excellence in Information Integrity Award Nomination*

**intelligent** enterprise REAL Awards **TRANSFORM**
**2004 WINNER**

*Best Business Intelligence Application Information Integration Client: Department of Defense*

**Chicago Tribune**

**Federal Computer Week**
GADGET

**DM Review**
2005

*World Class Solutions Award Data Management*

For more information on our Strategic Consulting Services, Implementation Services, or World-Class Training, call toll free at 866.EWS.1100, 866.397.1100, mai number 630.920.0005 or email us at Info@EWSolutions.com

www.EWSolutions.com

Strategic Partner & Systems Integrator

© 2008 Enterprise Warehousing Solutions, Inc. (EWSolutions) – 2

Intelligent Business Intelligence℠

## Professional Profile / Contact Information

**Pete Stiglich** is a Senior Consultant with EWSolutions with nearly 25 years of IT experience in the fields of Data Modeling, Data Warehousing, Business Intelligence, meta data Management, Data Integration, Customer Relationship Management (CRM), Customer Data Integration (CDI), Database Design and Administration, Data Quality, and Transaction Processing. Pete has architected Enterprise Information Management solutions for diverse industries such as Insurance, Credit Card, Medical, Retail, Banking, Manufacturing, Telecom, and Government.

Pete has developed and taught courses on Dimensional Data Modeling, Conceptual Data Modeling, ER/Studio, and SQL. Pete has presented for DAMA at the international and local level, as well as at the 2007 IADQ Conference. Pete's articles on Data Architecture have been published in *Real World Decision Support, DMForum, InfoAdvisors, and the Information and Data Quality Newsletter.* Pete is a listed expert in SearchDataManagement on the topics of data modeling and data warehousing.

For the current issue of Real World Decision Support

**See: http://www.ewsolutions.com/resource-center/rwds_folder/rwds-curr-issue/**

Email:     **PStiglich@EWSolutions.com**  Phone: **602-284-0992**

## EWSolutions' Partial Client List

| | | |
|---|---|---|
| Arizona Supreme Court | Ford Motor Company | Neighborhood Health Plan |
| Bank of Montreal | GlaxoSmithKline | NORC |
| BankUnited | Harris Bank | Physicians Mutual Insurance |
| Basic American Foods | The Hartford | Pillsbury |
| Becton, Dickinson and Company | Harvard Pilgrim HealthCare | Quintiles |
| Blue Cross Blue Shield companies | Health Care Services Corporation | Sallie Mae |
| Branch Banking & Trust (BB&T) | Hewitt Associates | Schneider National |
| British Petroleum (BP) | HP (Hewlett-Packard) | Secretary of Defense/Logistics |
| California DMV | Information Resources Inc. | South Orange County Community College |
| College Board | International Paper | SunTrust Bank |
| Corning Cable Systems | Janus Mutual Funds | Target Corporation |
| Countrywide Financial | Johnson Controls | The Regence Group |
| Defense Logistics Agency (DLA) | Key Bank | Thomson Multimedia (RCA) |
| Delta Dental | LiquidNet | United Health Group |
| Department of Defense (DoD) | Loyola Medical Center | United States Air Force |
| Driehaus Capital Management | Manulife Financial | United States Navy |
| Eli Lilly and Company | Mayo Clinic | United States Transportation Command |
| Federal Aviation Administration | Microsoft | USAA |
| Federal Bureau of Investigation (FBI) | National City Bank | Wells Fargo |
| Fidelity Information Services | Nationwide | Wisconsin Department of Transportation |
| | | Zurich Cantonal Bank |

**For more information on our Strategic Consulting Services, Implementation Services, or World-Class Training, call toll free at 866.EWS.1100, 866.397.1100, main number 630.920.0005 or email us at Info@EWSolutions.com**

The MIT 2008 Information Quality Industry Symposium

# What will we talk about?

- Data Models and Data/Information Quality

- What is a Conceptual Data Model?

- Benefits of Conceptual Data Models for Information Quality

- Developing the Conceptual Data Model

- Phased modeling approach (conceptual, logical, physical)

- Conceptual Data Model expressiveness

*www.EWSolutions.com*
© 2008 Enterprise Warehousing Solutions, Inc. (EWSolutions) – 5
*Strategic Partner & Systems Integrator*
Intelligent Business Intelligence℠

---

The MIT 2008 Information Quality Industry Symposium

# Data Models and Quality

*www.EWSolutions.com*
© 2008 Enterprise Warehousing Solutions, Inc. (EWSolutions) – 6
*Strategic Partner & Systems Integrator*
Intelligent Business Intelligence℠

The MIT 2008 Information Quality Industry Symposium

# Information and Data Quality

- Information and Data Quality is a <u>huge</u> issue for every business, government, or institution.

- Poor Information and Data Quality affects <u>every</u> type of information system – OLTP or decision support

- Often leads to a lack of confidence and credibility of IT and IT systems.

The MIT 2008 Information Quality Industry Symposium

# Information and Data Quality

- What is Data Quality?

  - Accurate, complete, and valid data that is captured, stored and maintained according to business requirements.

- What is Information Quality?

  - First, what is the difference between data and information?

**The MIT 2008 Information Quality Industry Symposium**

DATA ➕ META DATA ➕ DATA QUALITY ➕ DATA CONTEXT

| pst_cd |
|---|
| A7E 9U1 |
| B7I   6E2 |
| B6T 0X9 |
| 85016-0341 |
|  |
| DD8T 5V3 |

Policy Holder Postal Code

"The postal code of the policyholder.  Only Canadian postal codes in the format of "ANA NAN" is allowed.  Only 1 space is allowed between the 2 sections of the postal code.  (A represents an Alpha character, N represents a Numeric character).  If the Postal Code is unknown or invalid, the value "XXX XXX" will be substituted in order to facilitate filtering and data correction.

| plcyhldr_postal _cd |
|---|
| A7E 9U1 |
| B7I  6E2 |
| B6T 0X9 |
| XXX XXX |
| R3A A7E |
| D8T 5V3 |

Policy

Insured

Claim

= Information

---

**The MIT 2008 Information Quality Industry Symposium**

- Information Quality allows us to ask (and answer with confidence) questions such as?

  ➡ How many unique customers do we have across all lines of business?

  ➡ What geography would be the best to focus on for a new marketing campaign?

  ➡ What are patterns to look for in order to identify a potential disease outbreak?

  ➡ etc, etc, …

The MIT 2008 Information Quality Industry Symposium

# Information and Data Quality

■ There are many causes of poor <u>Data</u> Quality

➡ Lack of system constraints when data is originally captured

➡ Focus on quantity not quality (let's get these projects done as quickly as possible, and move on to the next thing…)

➡ Poor data management practices, e.g. authorization, archival

➡ Programmatic bugs

➡ Lack of management support for Data Governance and Stewardship

➡ Data Profiling tool not acquired/used!

➡ Lack of <u>automated</u> audits and alerts when actual/potential data quality events occur

➡ Etc….

---

The MIT 2008 Information Quality Industry Symposium

# Information and Data Quality

■ There are many causes of poor <u>Information</u> Quality

➡ Stovepiped, independent data marts – different people get different numbers for the same data

➡ Lack of an integrated Enterprise Data Warehouse, with <u>dependant</u> data marts

➡ Data not structured in an easy to use format (e.g. Dimensional) that can help prevent misunderstandings

➡ Users directly querying (e.g. via SQL tools) databases

➡ Lack of a Managed Meta Data Environment (MME)

   ■ What does this data mean?

   ■ Where did it originate from?

   ■ What were the conditions of the data at the time of the query – e.g. were any loads delayed

➡ Lack of Data Governance and Stewardship

➡ Etc…

The MIT 2008 Information Quality Industry Symposium

# Information and Data Quality

- However, an often overlooked cause of poor information and data quality is:

Poor or non-existing data models

Especially Conceptual Data Models!!

---

The MIT 2008 Information Quality Industry Symposium

# Data models and quality

- Data models are often an afterthought or developed only to meet immediate requirements.

- Data Models are often developed by application developers or DBA's – not by Data Architects.

- It is very common (and very bad practice) to see physical data models being the only data model developed for a system.  Better practice is to develop a logical model before a physical – but this is still not BEST practice!!

- Physical data models are optimized for performance – NOT for understandability.  Often, foreign key relationships are not utilized in Physical Data Models – making the physical model difficult to understand.

The MIT 2008 Information Quality Industry Symposium

# Data models and quality

- The physical data model, forward engineered to become the database schema may be in place for years or decades!!!

- Often much easier to change a program than to change a data model once a system is operational (or even while still in development)

- Ergo, data models should be developed with due rigor following industry best practices

  **Best Practice** is to use a phased modeling approach – conceptual, logical, and finally physical models

*www.EWSolutions.com*  
© 2008 Enterprise Warehousing Solutions, Inc. (EWSolutions) – 15  
*Strategic Partner & Systems Integrator*  
Intelligent Business Intelligence<sup>sm</sup>

---

The MIT 2008 Information Quality Industry Symposium

# Data models and quality

- What are some of the data and information quality issues that can arise from poor data models?

  - The application <u>does not meet business expectations</u>. Rework often required.

  - The model may meet the immediate needs of the application but may miss the larger needs of the enterprise.

  - M:M relationships may be missed which can lead to significant data duplication/missing data and increased development and maintenance costs

  - Business rules not identified, or not identified well. Business exceptions not identified possibly causing system outages.

  - If cardinality, optionality not properly identified, database constraints may be configured inaccurately leading to data quality problems.

  - If relationship identification not properly captured, granularity may be affected - data not being captured at the detail necessary, other problems.

  - Lack of good business meta data (attributes in business terms, business descriptions, identified data steward, etc)

  - More…

*www.EWSolutions.com*  
© 2008 Enterprise Warehousing Solutions, Inc. (EWSolutions) – 16  
*Strategic Partner & Systems Integrator*  
Intelligent Business Intelligence<sup>sm</sup>

The MIT 2008 Information Quality Industry Symposium

# Data models and quality

- A bad data architecture practice is developing Physical Data Models without developing Conceptual and Logical Data Models first

- IT needs to "Resist the Urge" to design physical (and logical) data models first.

---

The MIT 2008 Information Quality Industry Symposium

# Resisting the Urge

## What does this mean?

- There is a tendency to build physical data models first and ask questions later!!

- Not uncommon to see database schemas being developed in tandem with the application development process

- These models may meet initial requirements but break down when additional requirements and functionality are identified

- **These models often allow or even force Data Quality problems to creep in**

- Need to develop a conceptual data model as the first step of a phased modeling approach and use the conceptual data model as a tool to validate and communicate understanding of business requirements with the business

The MIT 2008 Information Quality Industry Symposium

# Causal factors

- Lack of data modeling experience and training

- IT professionals often don't feel productive unless they're "doing something" – e.g. developing a database or writing code.

- Temptation to cut corners when management wants things done yesterday

- Designing and creating databases is fun!!  Why did we get into IT but to design and build systems?

- In IT, there are many ways that something can be accomplished – not each way is equal in value

---

The MIT 2008 Information Quality Industry Symposium

# Result

- Systems which may not fully meet business requirements
- Physical structures that may initially be easy to load and query but over time become more difficult to use
- **Poor data quality!**
- Maintenance headaches
- Inflexible for future change
- Longer load cycles
- Etc…

**END RESULT:  Unsatisfied customers, increased expense, lack of confidence in IT, etc**

The MIT 2008 Information Quality Industry Symposium

**Group Exercise**

---

The MIT 2008 Information Quality Industry Symposium

## Example

### SCENARIO:

- A pet hospital chain that performs services and sells products started a CRM (Customer Relationship Management) undertaking and began capturing information about customers and their pets in a CDI (Customer Data Integration) Hub

- Also wanted to track household activity.  Last name and address used for determining a household.  A household is comprised of 1 or many customers.

- Data to be used for targeted marketing campaigns

- Wanted to be able to track multiple addresses per customer.

- Per business requirements, a Customer had only 1 household id

The MIT 2008 Information Quality Industry Symposium

## Example

**FACTORS:**

- Developers assumed they understood the business – they interviewed the customer

- A CDM was not created due many factors such as lack of data modeling expertise and tight deadlines.

- Was <u>incredibly</u> difficult to make changes to the model

- This "proof of concept" required very extensive modification in order for the business to have some confidence in it . It was eventually outsourced!

---

The MIT 2008 Information Quality Industry Symposium

## Example

**END RESULT:**

- Duplication all over the place, requiring unnecessarily complex processing and longer ETL processing windows

- Took heroic effort and a long amount time to adjust the system for changing business requirements – <u>CMM Level 0!</u> ↓

- Excessive maintenance programming

- <u>The business rules had to be enforced primarily in the ETL and SQL and not in the database!</u>

- The poor data model ***<u>forced</u>*** data quality problems into the system

- The data model didn't fulfill its "enforcement" role – enforcing good data quality through the data model!!

The MIT 2008 Information Quality Industry Symposium

# Headache

- As the old saying goes "An ounce of prevention is worth a pound of cure"

- Taking additional time up front to understand the business and develop **conceptual data models** helps:

  → Prevent <u>assumptions</u> which lead to data, information quality problems

  → Uncovers "gotchas" that can surface later – fewer "OH SHOOT" moments

  → Reduce development and maintenance costs

---

The MIT 2008 Information Quality Industry Symposium

# 7 Habits

- One of the habits in Steven Covey's **"7 Habits of Highly Effective People"** that is commonly quoted is *"Begin with the end in mind"*

- This makes great sense for many things but for good data modeling, start with the <u>beginning</u> in mind with an eye to the end (e.g. to limit scope for the CDM effort)

- *Understand the business first* and <u>finally</u> build physical structures (with many steps and iterations of steps in between)

- Understand the business first by developing a CDM, and review the CDM with the business

The MIT 2008 Information Quality Industry Symposium

# **What is the Conceptual Data Model?**

---

The MIT 2008 Information Quality Industry Symposium

## What is a Conceptual Data Model?

A diagram identifying real world concepts/objects/things (entities) and the relationships between these in order to gain, reflect, and document understanding of the business (as-is & to-be), in order to:

➡ foster semantic reconciliation

➡ improve business/IT collaboration

➡ serve as a framework for the development of information systems

The MIT 2008 Information Quality Industry Symposium

# What is a Conceptual Data Model?

"*A conceptual entity-relationship model shows how the business world sees information. It suppresses non-critical details in order to emphasize business rules and user objects. It typically includes only significant **entities** which have business meaning, along with their **relationships**.* "

**Applied Information Science website**

---

The MIT 2008 Information Quality Industry Symposium

# What is a Conceptual Data Model

'**A data model that represents an abstract view of the real world. A conceptual model represents the human understanding of a system…. A conceptual data model describes how relevant information is structured in the natural world. In other words, it is how the human mind is accustomed to thinking of the information.**'

**OECD Glossary of Statistical Terms**

The MIT 2008 Information Quality Industry Symposium

# What is a Conceptual Data Model?

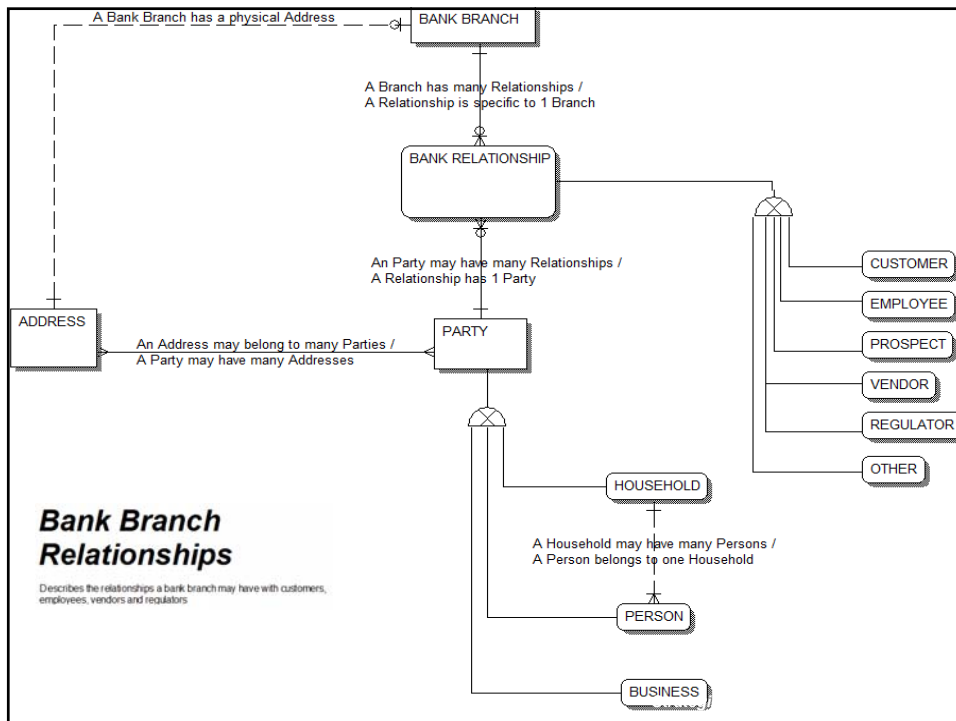- It is "stateless" - <u>NOT</u> a state model

- The <span style="color:red">entire</span> <span style="color:teal">possible</span> <span style="color:green">lifecycle</span> of a relationship should be represented, <u>per current business practice</u>

➡ This includes <u>business</u> <span style="color:red">exceptions!!</span>

➡ *Not exceptions due to poor data quality or due to system limitations*)

➡ The CDM should reflect the <u>business</u> – not IT systems

➡ Review optionality and cardinality to ensure longitudinal perspective

*www.EWSolutions.com*
© 2008 Enterprise Warehousing Solutions, Inc. (EWSolutions) – 33

*Strategic Partner & Systems Integrator*
Intelligent Business Intelligence℠



A Bank Branch has a physical Address — BANK BRANCH

A Branch has many Relationships /
A Relationship is specific to 1 Branch

BANK RELATIONSHIP

An Party may have many Relationships /
A Relationship has 1 Party

ADDRESS

An Address may belong to many Parties /
A Party may have many Addresses

PARTY

CUSTOMER
EMPLOYEE
PROSPECT
VENDOR
REGULATOR
OTHER

HOUSEHOLD

A Household may have many Persons /
A Person belongs to one Household

PERSON

BUSINESS

**Bank Branch Relationships**

Describes the relationships a bank branch may have with customers, employees, vendors and regulators

The MIT 2008 Information Quality Industry Symposium

# Semantic Resolution

- A CDM is a key tool for semantic resolution

- For enterprise applications, have to reach consensus across divisions, departments, external agencies, etc, for naming and defining data entities, and identifying correct relationships.

- Semantic resolution is a key activity of Data Governance and Stewardship, and an ECDM is a key enabler of Data Governance and Stewardship – these activities often take place in tandem, iteratively

- **Difficult to have Information Quality if synonyms, homonyms haven't been resolved.** *E.g. Is a customer a party that has placed an order, or can customer be a party who placed an order or a party that might become a paying customer?*

---

The MIT 2008 Information Quality Industry Symposium

# Semantic Resolution

- Due to fundamental differences with the LDM, the CDM often has to be contained in a separate model file and **so there is a risk that lineage from a logical entity to a conceptual entity can be lost**

- Be sure to save the association between conceptual and logical entities, logical and physical entities, etc using:

  - A meta data repository and related tool which can be used to establish these relationships
  - User defined meta data properties within the model
  - Spreadsheet, etc. Last resort

- CDM's can help drive creation of a common, corporate lexicon – fostering improved communication, standardization --- BENEFICIAL TO THE ENTIRE ENTERPRISE – NOT JUST IT!

Gary Larson – The Far Side

The MIT 2008 Information Quality Industry Symposium

# Developing the Conceptual Data Model

---

The MIT 2008 Information Quality Industry Symposium

## Getting started developing a CDM

- A major hurdle is separating "data thinking" vs "process thinking"

- For conceptual data modeling, we're thinking about **"what"** (data) not the **"how"** (process).

- For a CDM – data is a relative term

- Data may not exist currently for a conceptual entity – but entities must be included in the CDM if it is an object of importance to the business

*This – not this*

THING1    THING2

The MIT 2008 Information Quality Industry Symposium

# Getting started developing a CDM

- When interviewing the business helpful to use a "recipe" analogy (see Steve Hoberman design challenge *) . A recipe identifies the **ingredients, utensils, equipment (whats)** and has **directions (hows)** in order to meet the desired goal.

- If the interviewee focuses on process ask "What things are needed for the XYZ process?" "What are the components of the XYZ process?"

- Helpful starting place is to identify "nouns", e.g. Customer, Product, Inventory

 **\*** DMReview January 2008, quoting Geof Clark

---

The MIT 2008 Information Quality Industry Symposium

# When is a CDM finished?

**"Perfection does not come into being, when nothing more can be added, but when nothing can be taken away"**

Antoine de Saint-Exupéry

The MIT 2008 Information Quality Industry Symposium

# Information Engineering (IE) notation

*Entity*

*Optionality*

*Relationship Identifying or Non-Identifying*

*Entity Name*

Organization

System

*Mandatory or Optional*

Organization ID

Utilizes

System ID

*Attributes*

Organization Name

Organization ID

*Primary Key*

Taxpayer ID (AK)

System Name

*One*

*Many*

*Verb Phrase*

*Alternate Key*

*Cardinality*

---

The MIT 2008 Information Quality Industry Symposium

# Entity (type)

- In the CDM, a real-world object of interest to the business

- Don't include associative entities or entities that mirror a database table (*unless it corresponds to a real-world business object*)

**Identifying Relationship**

**Independent**

**(square)**

**Dependant**

ORDER

ORDER NUMBER

ORDER LINE

ORDER NUMBER (FK)
ORDER LINE NUMBER

**(rounded edges)**

The MIT 2008 Information Quality Industry Symposium

# Identification in Relationships

──────── **An identifying relationship is stronger -- helps determine the meaning and granularity of a child entity. Is always mandatory. (NOTE: a solid line in a M:M relationship does not denote an identifying relationship!!)**

── ── ── **A non identifying relationship may be mandatory or optional, but does not define meaning/granularity**

Identifying Relationship        Non-identifying Relationship

ORDER_HEADER
🔑 ORDER_NUMBER

ORDER_LINE
🔑 ORDER_NUMBER (FK)
🔑 ORDER_LINE_NUMBER

ORDER TYPE
ORDER TYPE ID

ORDER HEADER
ORDER NUMBER
ORDER TYPE ID (FK)

*www.EWSolutions.com*
© 2008 Enterprise Warehousing Solutions, Inc. (EWSolutions) – 47

*Strategic Partner & Systems Integrator*
Intelligent Business Intelligence$^{sm}$

---

The MIT 2008 Information Quality Industry Symposium

# Relationship Verb Phrase

- Describes the relationship using business terminology

- Can be terse but IMO verbose is better

- <u>You never know who will end up looking at the model!!!</u>

- Business people probably won't take the time to understand the notation!

BANK

A Bank has many Accounts /
An Account belongs to one Bank

ACCOUNT

An Account may have many Customers /
A Customer may have many Accounts

CUSTOMER

*www.EWSolutions.com*
© 2008 Enterprise Warehousing Solutions, Inc. (EWSolutions) – 48

*Strategic Partner & Systems Integrator*
Intelligent Business Intelligence$^{sm}$
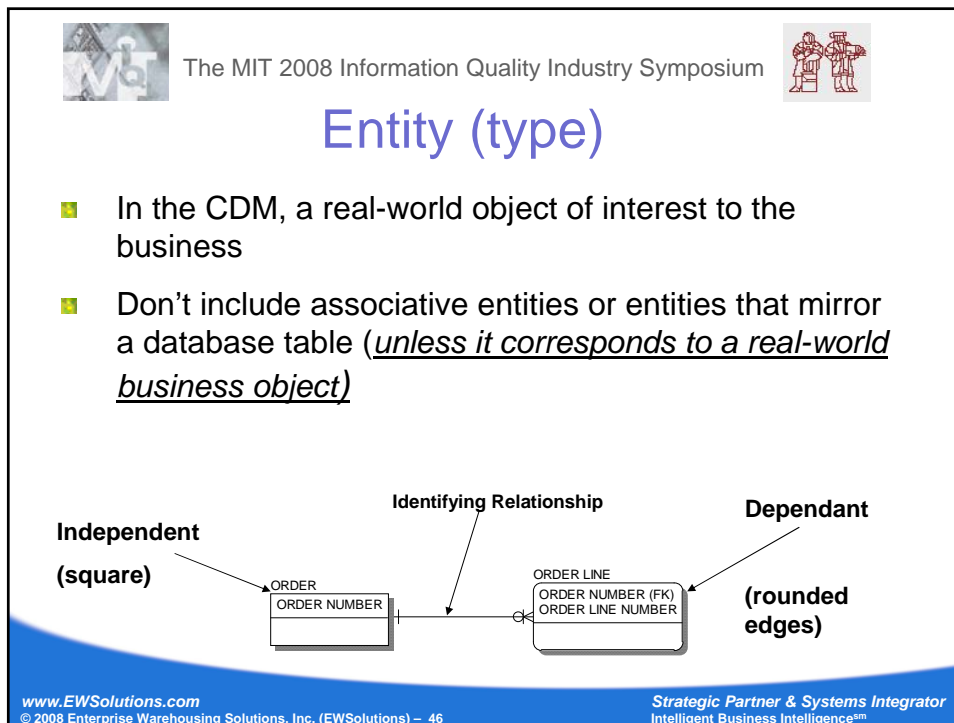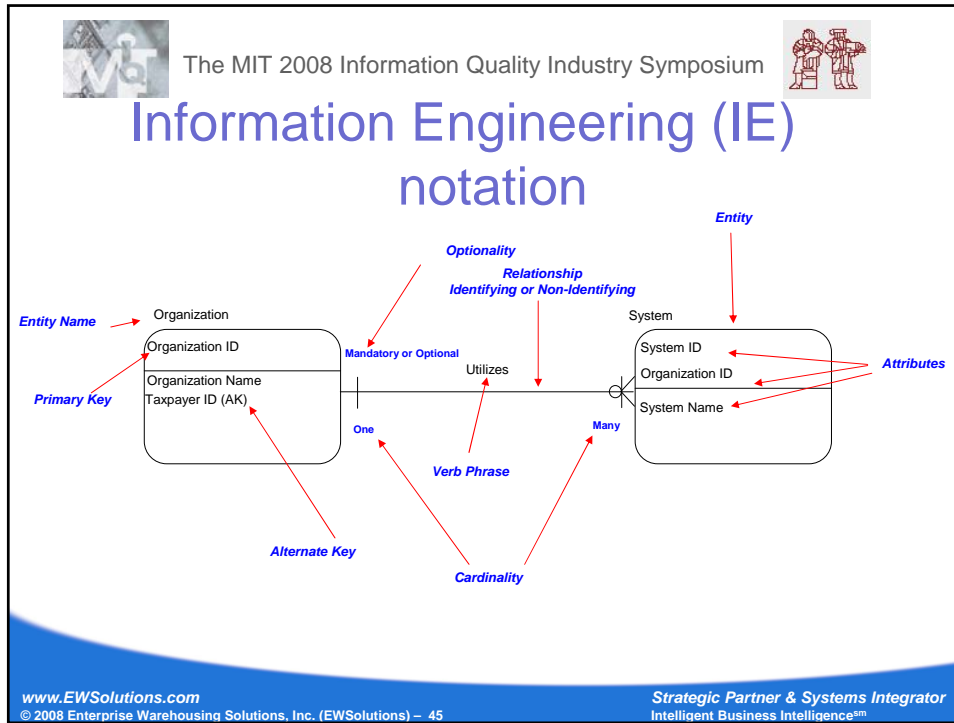
The MIT 2008 Information Quality Industry Symposium

# Subtyping

- Subtypes/Supertypes make a model more <u>expressive and understandable</u>

- Subtypes describe a Supertype

- A Subtype can be inclusive or exclusive, exhaustive or non-exhaustive

CUSTOMER

Exclusive Subtype

PROSPECTIVE CUSTOMER

CUSTOMER WITH ACCOUNT

CUSTOMER

Inclusive Subtype

LOAN CUSTOMER

DEPOSIT ACCOUNT CUSTOMER

**More on subtyping later**

---

The MIT 2008 Information Quality Industry Symposium

# Relationships and quality

Many to Many Relationships

- A conceptual data model will <u>very often</u> have numerous M:M relationships in order to accurately reflect all possible states of a relationship

- A CDM is not a state model – it should reflect the relationship from a longitudinal (entire lifecycle of the relationship) perspective

- For example, a store clerk works for one store in almost all cases, but it is possible for a clerk to move and begin work with another store.

Clerk:Store  s/b a M:M

The MIT 2008 Information Quality Industry Symposium

# Relationships and quality

Many to Many Relationships

- When a M:M relationship is not identified during requirements definition in a CDM….

    - Project scope is not measured correctly

    - Logical model design, application development, testing are all impacted – <u>heavily</u>!!!

    - Some M:M instances occur only occasionally – can cause bugs, outages, missed or duplicate data weeks/months later when exceptions are encountered
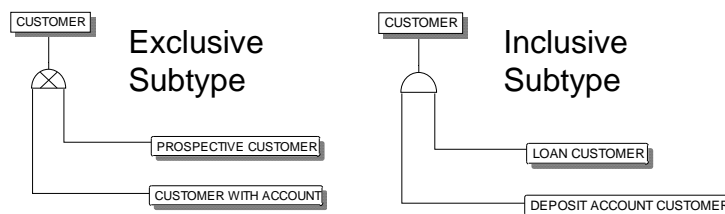
*www.EWSolutions.com*
© 2008 Enterprise Warehousing Solutions, Inc. (EWSolutions) – 51

*Strategic Partner & Systems Integrator*
Intelligent Business Intelligence<sup>sm</sup>

---

The MIT 2008 Information Quality Industry Symposium

# Relationships and quality

Many to Many Relationships

- Impacts

    - Logical and physical models have to be revisited, reviewed, and possibly reapproved
    - Can have a tremendous impact on applications – screen forms, program functions, load processes, reports, SQL, cubes, etc..
    - Existing data may need to be restructured
    - Impacts to downstream systems (DW/BI, ODS, MDM, etc)

*www.EWSolutions.com*
© 2008 Enterprise Warehousing Solutions, Inc. (EWSolutions) – 52

*Strategic Partner & Systems Integrator*
Intelligent Business Intelligence<sup>sm</sup>

The MIT 2008 Information Quality Industry Symposium

# Relationships and quality

- Resolving M:M relationships involves non-trivial decisions, with benefits and impacts to weigh. Not something to decide during a 3 am support call…..

- Resolution decision can have a dramatic impact on quality

- If you choose not to allow M:M relationship in a particular instance – how are you going to impact the business near or long term?

---

The MIT 2008 Information Quality Industry Symposium

# Relationships and quality

- A <u>major</u> bank allowed multiple customers to apply for a single loan

  - ➡ However only the information for the 1[st] customer was <u>retained</u> (e.g. identifying information, credit score) in the system.

  - ➡ Bank had **no <u>accurate</u> idea how many customers it had**, and could not easily and accurately gauge the total customer experience…..

  - ➡ Might not know if it was marketing a new loan to a customer who had defaulted on a prior loan…

The MIT 2008 Information Quality Industry Symposium

# Relationships and quality

- Model relationships are a key source for of business rules and data quality metrics.

- For non-kernel entities, **identifying relationships** will be critical to understanding entity meaning and granularity – be sure to distinguish identifying / non-identifying

- **CAN HAVE A DRAMATIC IMPACT ON UNDERSTANDING (or misunderstanding) THE MODEL!!**

---

The MIT 2008 Information Quality Industry Symposium

# Relationships and quality

- If you want your application to be successful…

## Data relationships must be correctly identified!!!!

- The only question is: <u>when are you going to pay</u> to discover the correct relationship???

The MIT 2008 Information Quality Industry Symposium

# Relationships and quality

- You can identify the <u>correct</u> relationships

$ In the CDM phase (during requirements definition)

$ $ During Logical Data Model development

$ $ $ During Physical Data Model development

$ $ $ $ During application development

$ $ $ $ $ During implementation

$ $ $ $ $ $ $ $ $ During production

The MIT 2008 Information Quality Industry Symposium

# CDM and requirements

- Some statistics…

- If it costs $1 to fix a defect found in the requirements phase, it costs $2 in design phase, and continues to rapidly increase until it costs $68 if not found until product is released into operation - *Boehm, Barry W. Software Engineering Economics. Englewood Cliffs, NJ: Prentice-Hall, 1981*

- Requirements errors account for 70 to 85 percent of the rework cost - *Leffingwell 1997, quoting Barry W. Boehm*

- The cost to fix the defect in QA stage is eight times more than during the Requirements Development stage - *Grady 1999*

The MIT 2008 Information Quality Industry Symposium

# CDM and requirements

- Conceptual data modeling should take place in the <u>Requirements Definition</u> phase
- Conceptual data modeling (in general) is NOT <u>design</u> – it is <u>description</u>
- Modeling the <u>BUSINESS</u> – not a SOLUTION

The MIT 2008 Information Quality Industry Symposium

# Data Modeling Progression

The MIT 2008 Information Quality Industry Symposium

# Model progression

- Conceptual Data Model (CDM)

  - **Technology and application neutral**

  - Entities may or may not eventually translate into a physical database table

  - **A data source for a conceptual entity does not need to exist!!**  Only interested in understanding the business at this point

  - **Physical implementation is NOT important at this point** – conceptual data modeling is all about documenting business objects.  Set expectations appropriately when presenting to technologies personnel.

The MIT 2008 Information Quality Industry Symposium

# Model progression

- If the CDM is wrong, your downstream models may be built upon incorrect premises!!!

- Don't shortchange the amount of time spent in this step!!

- NEVER a waste of time!!  *At the very least you can justify it as a tool for yourself for developing LDM's – who can remember all the identification, cardinality, optionality of even a moderately complex subject area?*
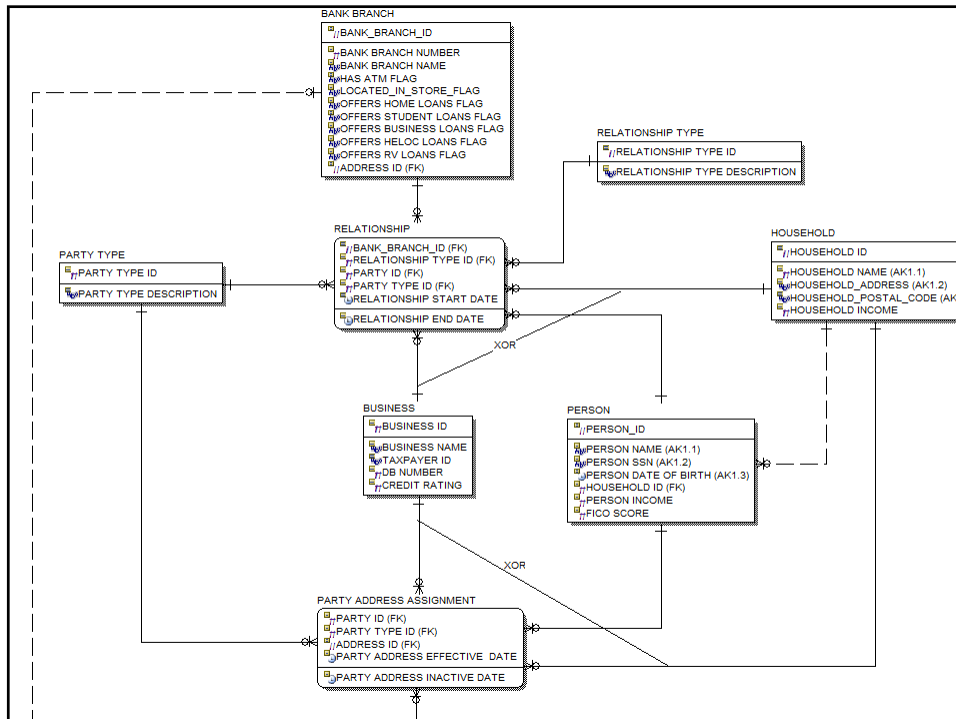
The MIT 2008 Information Quality Industry Symposium

# Model progression

- Logical Data Model (LDM)

  - First step of <u>SOLUTION DATA  DESIGN</u> (generally)

  - Fully/mostly attributizes a conceptual data model

  - Resolves many to many relationships (usually)

  - Resolves subtypes/supertypes (usually)

  - May introduce abstraction (generalize entities, attributes, relationships) – more later

  - Formalize primary keys

**BANK BRANCH**
- // BANK_BRANCH_ID
- BANK BRANCH NUMBER
- BANK BRANCH NAME
- HAS ATM FLAG
- LOCATED_IN_STORE_FLAG
- OFFERS HOME LOANS FLAG
- OFFERS STUDENT LOANS FLAG
- OFFERS BUSINESS LOANS FLAG
- OFFERS HELOC LOANS FLAG
- OFFERS RV LOANS FLAG
- // ADDRESS ID (FK)

**RELATIONSHIP TYPE**
- // RELATIONSHIP TYPE ID
- RELATIONSHIP TYPE DESCRIPTION

**HOUSEHOLD**
- // HOUSEHOLD ID
- HOUSEHOLD NAME (AK1.1)
- HOUSEHOLD_ADDRESS (AK1.2)
- HOUSEHOLD_POSTAL_CODE (AK1
- // HOUSEHOLD INCOME

**PARTY TYPE**
- // PARTY TYPE ID
- PARTY TYPE DESCRIPTION

**RELATIONSHIP**
- // BANK_BRANCH_ID (FK)
- // RELATIONSHIP TYPE ID (FK)
- // PARTY ID (FK)
- // PARTY TYPE ID (FK)
- RELATIONSHIP START DATE
- RELATIONSHIP END DATE

XOR

**BUSINESS**
- // BUSINESS ID
- BUSINESS NAME
- TAXPAYER ID
- // DB NUMBER
- // CREDIT RATING

**PERSON**
- // PERSON_ID
- PERSON NAME (AK1.1)
- PERSON SSN (AK1.2)
- PERSON DATE OF BIRTH (AK1.3)
- // HOUSEHOLD ID (FK)
- // PERSON INCOME
- // FICO SCORE

XOR

**PARTY ADDRESS ASSIGNMENT**
- // PARTY ID (FK)
- // PARTY TYPE ID (FK)
- // ADDRESS ID (FK)
- PARTY ADDRESS EFFECTIVE  DATE
- PARTY ADDRESS INACTIVE DATE

The MIT 2008 Information Quality Industry Symposium

# Abstraction in the CDM

- *"Abstraction is the removal of details in such a way as to broaden applicability to a wide class of situations while preserving the important properties and essential nature from concepts or subjects"  \**

- In the CDM, generally avoid abstraction in order to more closely mirror the business.

- Use supertypes when you need to abstract – for establishing broad applicability relationships (in order to avoid establishing relationships to all the subtypes)

\* *Steve Hoberman – Data Modeling Made Simple*

---

The MIT 2008 Information Quality Industry Symposium

# Abstraction in the CDM

**With abstraction..**          **Without….**



**Now add 20+ more types of parties (e.g. insurance) more relationships, ….**

The MIT 2008 Information Quality Industry Symposium

# Abstraction in the LDM

- In the LDM, abstraction is necessary for normalization – data stored only once

- Entities, attributes, relationships can be abstracted

- Allows for flexibility in case other types need to be added in the future

---

The MIT 2008 Information Quality Industry Symposium

# Model Progression

- **Physical Data Model (PDM)**

  - Represents how a logical model is applied to a particular DBMS platform

  - Assign datatypes, indexing, storage, partitioning, etc

  - Can be forward engineered to create the actual database structures

  - Complies with DBMS nomenclature restrictions

  - PDM may look different than the logical – e.g. column ordering to take advantage of partition elimination

The MIT 2008 Information Quality Industry Symposium

# Model Progression

- Why develop all these models?

  - Follows the progression in which a Data Modeling project should be undertaken

  - As more information becomes known, the more depth the models will be able to convey

  - Data Models convey knowledge – and knowledge is retained in data models
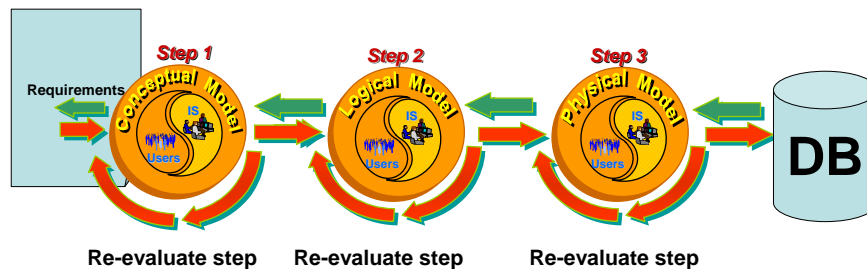
The MIT 2008 Information Quality Industry Symposium

# Phased Modeling Approach

- *Defined within the scope of the business problem*

- Data Modeling is an iterative endeavor – and it will probably be necessary to make revision to upstream models as more information becomes available

The MIT 2008 Information Quality Industry Symposium

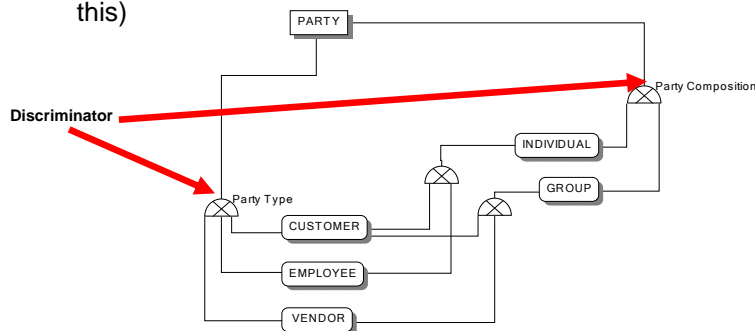# Conceptual Data Model Expressiveness

---

The MIT 2008 Information Quality Industry Symposium

## Subtyping

- An entity may have many subtype relationships - use a discriminator to distinguish the subtype relationship

- A Subtype may have more than one Supertype (not every tool allows this)

The MIT 2008 Information Quality Industry Symposium

# Subtyping

- **A single Subtype entity may have relationships with entities which do not apply to the other Subtypes**
- This helps the CDM to better mirror business reality.
- **<u>Additional business rules can be expressed!!</u>**
- When entity abstraction occurs during the LDM phase – entities and relationships are "lost". The identification, cardinality, and optionality of these relationships is subsumed into the remaining relationships
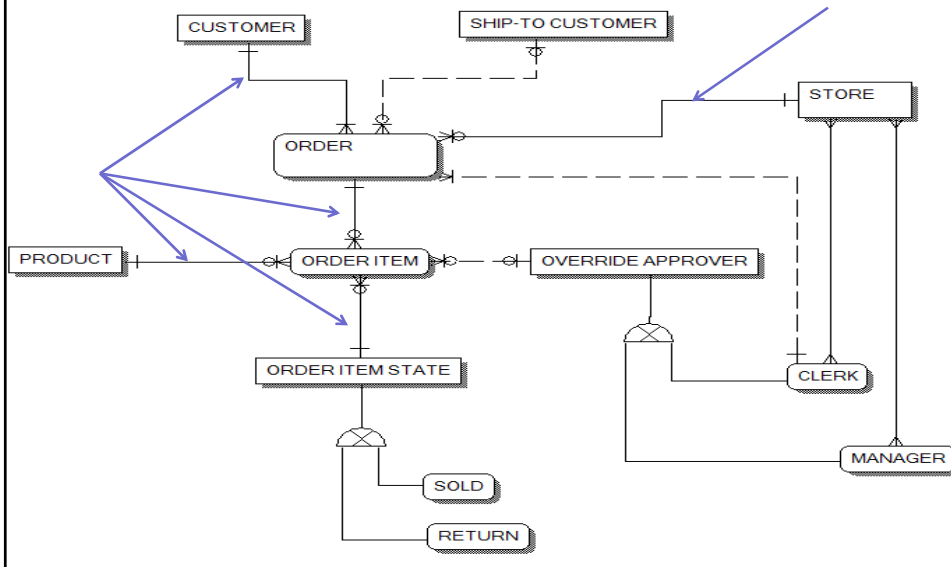
**NOTE: abstraction is a very useful and valid tool for the logical design phase - critical for normalization to eliminate redundant data. However, business semantics and rules are harder to identify – especially by a business person.**
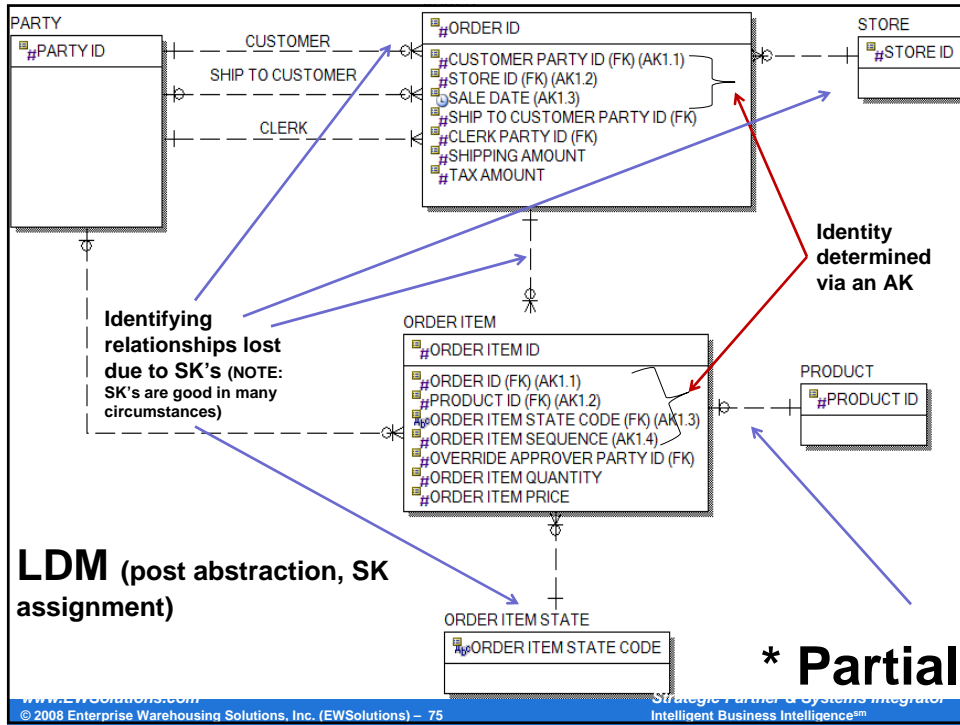
*www.EWSolutions.com*
© 2008 Enterprise Warehousing Solutions, Inc. (EWSolutions) – 73

*Strategic Partner & Systems Integrator*
Intelligent Business Intelligence℠

---

The MIT 2008 Information Quality Industry Symposium

# CDM (pre abstraction – note identifying relationships)

**LDM** (post abstraction, SK assignment)



# CDM Expressiveness

The MIT 2008 Information Quality Industry Symposium

- The previous model, which was a normalized logical data model, partially based off of the earlier CDM , is a solution model – not a business model.

- Identifying relationships are lost due to the surrogate key assigned to the "sale" entity

- Business rules (relationships) are still there, but aren't as obvious – it isn't modeled how "the human mind is accustomed of thinking about information"

The MIT 2008 Information Quality Industry Symposium

## Subtyping and Taxonomies

- Establishing subtypes in a CDM is often the first step in developing taxonomies for classifying data

- Provides value domain for a taxon (full or partial)

- Makes abstract names more understandable

**Framed subtyping (Euler diagram)**

**PRODUCT CATEGORY**
- TRUCK
- CAR
- SPORT UTILITY VEHICLE

**PRODUCT MAKE**
- CHRYSLER
- DODGE
- HONDA
- OTHER

**PRODUCT MODEL**
- PILOT
- OTHER
- SEBRING
- RIDGELINE

*www.EWSolutions.com*
© 2008 Enterprise Warehousing Solutions, Inc. (EWSolutions) – 77
*Strategic Partner & Systems Integrator*
Intelligent Business Intelligence[sm]

---

The MIT 2008 Information Quality Industry Symposium

## Wrapping it up

- There are many causes of poor Data and Information Quality

- Poor or non-existing CDM's are not the least of these causes

- Need to "Resist the urge" to develop physical (and logical) models before developing conceptual models

- Many business requirements and rules are captured and documented in the CDM

- The CDM is a key means to validate IT's understanding of business requirements, and can be used to measure data quality in implemented systems

- CDM's need to be validated by the business

- CDM's need to be presentable, understandable and tailored to the audience

*www.EWSolutions.com*
© 2008 Enterprise Warehousing Solutions, Inc. (EWSolutions) – 78
*Strategic Partner & Systems Integrator*
Intelligent Business Intelligence[sm]

The MIT 2008 Information Quality Industry Symposium

# Wrapping it up

**Developing CDM's <u>first</u> is beneficial to your organization!!**

- Models downstream from CDM's more accurately reflect business requirements

- Fosters semantic resolution, in turn improving Information Quality

- Relationship identification, cardinality, and optionality are critical to good Data Quality

- Development and maintenance work is simplified and costs are reduced. **Once a system goes into production, it is very hard to change data structures!**

*www.EWSolutions.com*
© 2008 Enterprise Warehousing Solutions, Inc. (EWSolutions) – 79
*Strategic Partner & Systems Integrator*
Intelligent Business Intelligence℠

---

The MIT 2008 Information Quality Industry Symposium

# References

- Applied Information Science website

- OECD Glossary of Statistical Terms

- Zachman Framework for Enterprise Architecture

- Steve Hoberman Design Challenge, DMReview January 2008, quoting Geof Clark

- Boehm, Barry W. Software Engineering Economics. Englewood Cliffs, NJ: Prentice-Hall, 1981

- Steve Hoberman, "Data Modeling Made Simple" Technics Publications, 2005

*www.EWSolutions.com*
© 2008 Enterprise Warehousing Solutions, Inc. (EWSolutions) – 80
*Strategic Partner & Systems Integrator*
Intelligent Business Intelligence℠

The MIT 2008 Information Quality Industry Symposium
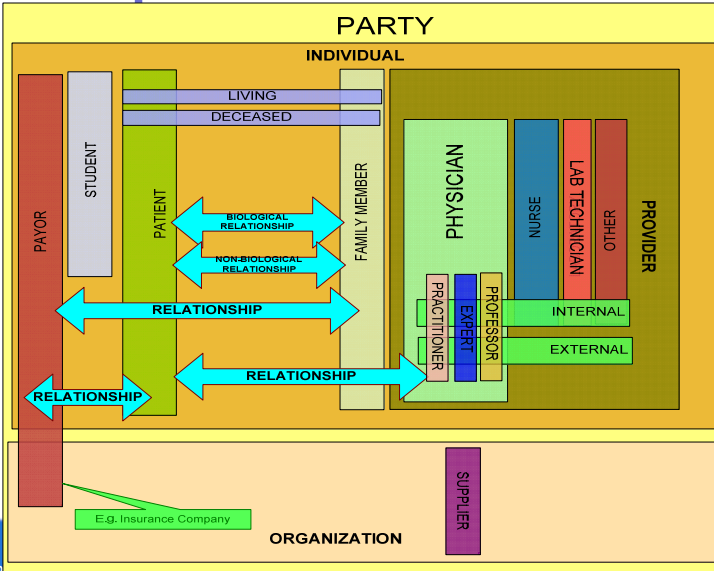
**Questions?**

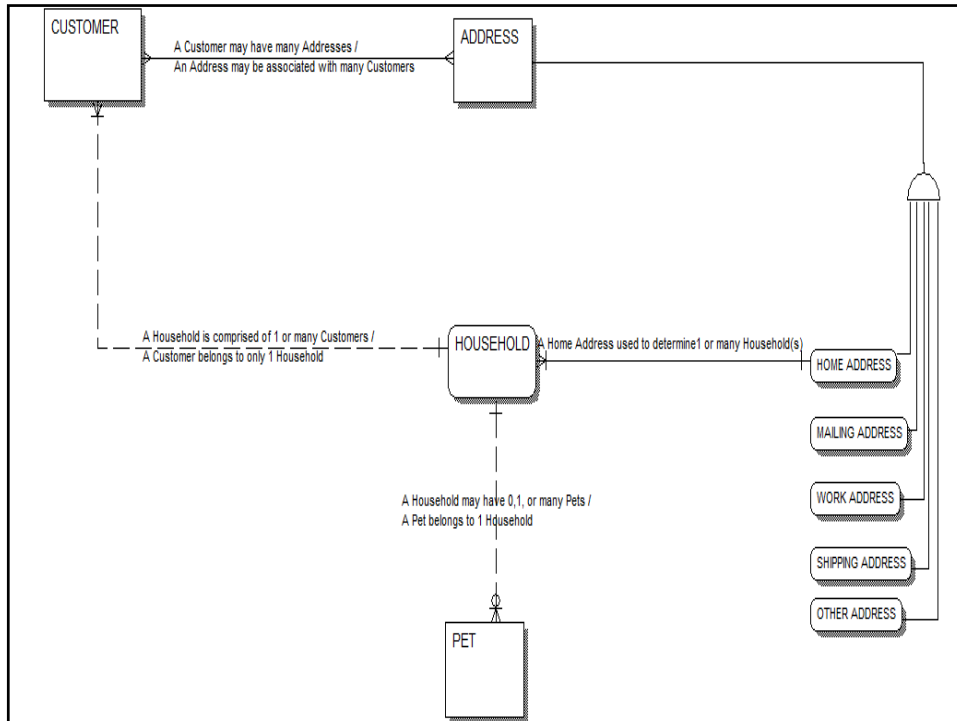The MIT 2008 Information Quality Industry Symposium

**Graphic Model**

Sample graphical representation of a CDM regarding the Party subject area in a medical educational institution

**Which are inclusive subtypes and which are exclusive subtypes?**

PARTY

INDIVIDUAL

PAYOR

STUDENT

PATIENT

LIVING

DECEASED

FAMILY MEMBER

PHYSICIAN

NURSE

LAB TECHNICIAN

OTHER

PROVIDER

PRACTITIONER

EXPERT

PROFESSOR

INTERNAL

EXTERNAL

BIOLOGICAL RELATIONSHIP

NON-BIOLOGICAL RELATIONSHIP

RELATIONSHIP

RELATIONSHIP

RELATIONSHIP

SUPPLIER

E.g. Insurance Company

ORGANIZATION

The MIT 2008 Information Quality Industry Symposium

**EWSolutions, Inc.**
**15 Spinning Wheel Road,**
**Suite 330**
**Hinsdale, IL 60521**
**Office   630.920.0005**
**Fax      630.920.0008**

**http://www.EWSolutions.com**