

GE DAMA-Phoenix, 2009/09

35START
1

Subtypes & Supertypes

The Data Modeler's most Valuable Construct



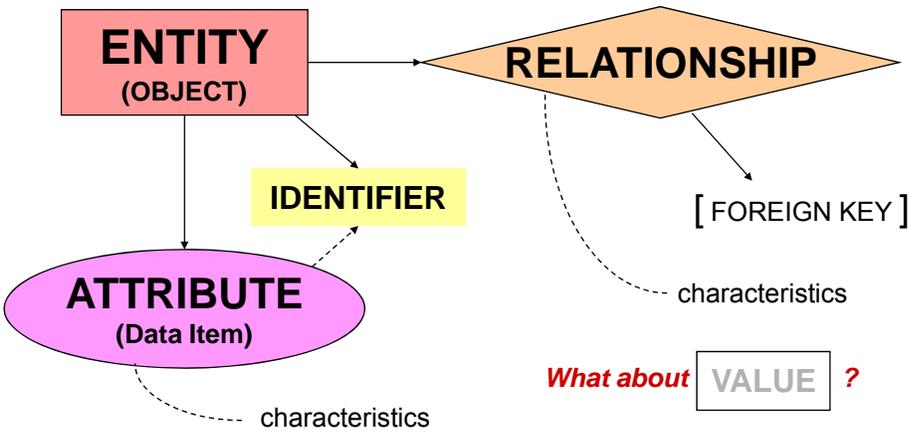
©Gordon C. Everest
Professor Emeritus of MIS and Database
Carlson School of Management
University of Minnesota
geverest@umn.edu

© Gordon C. Everest, All rights reserved.

GE Data Modeling Constructs

DMOD
2

What to look for.
Relative emphasis differentiates Data Modeling Schemes
e.g. ER modeling focuses on Entities and Relationships,
de-emphasizing, even hiding Attributes.



The diagram illustrates the relationships between data modeling constructs:

- ENTITY (OBJECT)** (red rectangle) connects to **RELATIONSHIP** (orange diamond) and **ATTRIBUTE (Data Item)** (pink oval).
- RELATIONSHIP** connects to **IDENTIFIER** (yellow rectangle) and **[FOREIGN KEY]** (text).
- ATTRIBUTE (Data Item)** connects to **IDENTIFIER** and **VALUE** (grey rectangle).
- Dashed lines labeled "characteristics" point from **RELATIONSHIP** and **ATTRIBUTE (Data Item)** to **VALUE**.
- A red text prompt asks: *What about* **VALUE** ?

 <small>431START</small>	Advanced Database Design
3	<h2 style="text-align: center;">7. Sub/Super Types</h2> <ul style="list-style-type: none">• “Abstractions” v. Collections v. Generalization• Attribute, Relationship, & Entity Generalization• Subtypes and Supertypes (Entity Generalization)<ul style="list-style-type: none">– Underlying assumptions; conditions– Generalization vs. Specialization– Diagramming - Graphical representations• Constraints• Subtype Definition - distinguishing attribute• Sub/SuperType Hierarchy/Lattice<ul style="list-style-type: none">– The Universal Relation• Inheritance (single; multiple) & Reuse• Mapping to Tables <hr/> <p style="text-align: center;">©Gordon C. Everest Carlson School of Management University of Minnesota</p> <p style="text-align: left;"><small>© Gordon C. Everest, All rights reserved.</small></p>

	Starting from Basics	
4	When you see this in an ER/schema diagram:	
		
	<p style="text-align: center;"><i>What does it mean?</i></p> <p style="text-align: center;"><i>What can you assume?</i></p> <p style="text-align: center;"><i>What does "the system" assume?</i></p>	
N		At least four distinct (context-free) semantic statements.



A Model

DMOD

5

A Model is an abstract re.presentation

... of something observed or something to be built.

Some "reality" (in some domain of interest) *presents* itself to the modeler.

The modeler builds a model of what they perceive or design.

Hence, the model is a *re.presentation* of that reality.

Abstract because some aspects of the "reality" are omitted

(necessarily, because people are doing the observing or designing and therefore it is an imperfect process).

The *semantics* of a model are seen through some *syntax*.

The *syntax* is the chosen method/scheme for representation.



Abstraction

DMODPRE

6

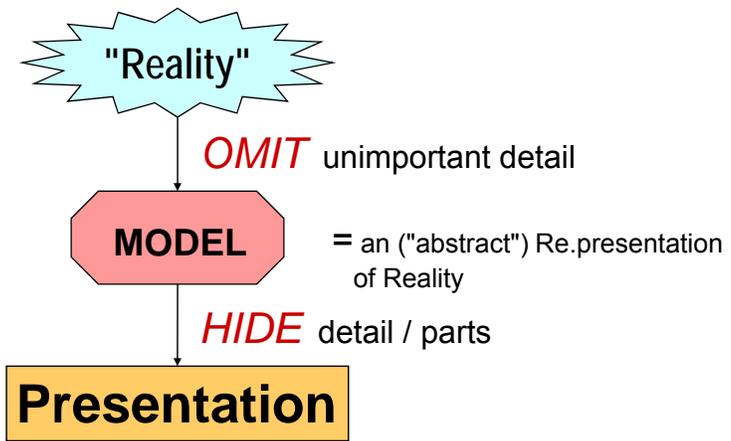
Jeff Kramer, CACM (50:4) 2007 April, p.36
www.cse.msu.edu/~cse914/Papers/Kramer-abstraction.pdf

- Definition - used in two senses:
 - (1) removing detail
 - (2) generalization - to identify commonalities
- Fit for a Purpose
- Widely used in art, music (jazz), maps (London Underground)
- Essential to Computing, (Software) Engineering, Requirements elicitation, (Data) Modeling
- Fourth level of human development (J. Piaget, 1896-1980)
 - "formal operational"- ability to think abstractly, systematically, hypothetically, symbolically
 - only one third achieve this level; some people never reach it!
They do not appreciate the value of modeling; and find it difficult to identify what is important in a problem.
- Need a test to measure abstraction ability



Abstraction: Omitting vs. Hiding

7 DMODPRE
SSTYPE



"There is no abstract art. You must always start with something. Afterward you can remove all traces of reality." -- Pablo Picasso



"Abstractions" & Collections

8 SSTYPE

Focusing on *selected* properties of objects

- Batini, Ceri, & Navathe
 - Smith & Smith (1977)
 - Graeme Simsion
 - Steve Hoberman
 - Len Silverston
 - David Hay

"Abstractions": (used in a different sense here)

- CLASSIFICATION ("Member-of")
 - Forming Types - entity sets/populations, domains
- AGGREGATION ("Part-of")
 - Building an entity record with descriptors (clustering attributes)
 - COMPOSITION (stronger "Part of" - no independent existence)
- GENERALIZATION/SPECIALIZATION ("Is-a")
 - Forming subtypes/supertypes, population subsets

Collections: (assumes *homogeneous* members)

- SET – no duplicates and no order (the only one in Relational)
- BAG – counting duplicates
- SEQUENCE – order matters
- ...

and looking at Relationships/Structure: GRAPH, TREE, ...

	<h2>Generalization</h2>	
<small>SSTYPE</small>		<small>see: Graeme Simsion, Ch. 4</small>
9	<ul style="list-style-type: none"> • Recognizing commonalities <ul style="list-style-type: none"> + valued - Do I care? Is it useful? - cost - Is it worth the effort? • Moving "up" to a higher, more inclusive, more generic, more "abstract" view <p>TYPES:</p> <ul style="list-style-type: none"> • Attribute <ul style="list-style-type: none"> - constrained by Entity Generalization - often a prelude to Entity Generalization • Entity <ul style="list-style-type: none"> - represented using subtypes/supertypes - implications for placement and naming of attributes and relationships • Relationship 	

	<h2>Attribute Generalization Examples</h2>	
<small>SSTYPE</small>		<small>Steve Hoberman, Data Modeler's Workbench</small>
10	<ul style="list-style-type: none"> • For a Tuxedo Rental shop, store Customer attributes: <ul style="list-style-type: none"> - Waist size - Leg length - Neck size - Arm length - Shoulder width <p>=> Later add Shoe size.</p> <p><i>What does that do to your database schema?</i></p> <p><i>How might you solve the problem?</i></p> <p><i>How does referential integrity become important here?</i></p> <p><i>What is the down side of this schema redesign?</i></p> <p>Similarly for Phone numbers:</p> <p>Problems:</p> <ul style="list-style-type: none"> - Handling international numbers - Handling other contact information, e.g. email 	
N		



Attribute Generalization Example

SSTYPE
Steve Hoberman, *Data Modeler's Workbench*

11 Given the following three entity type populations:

What do you observe?

CUSTOMER	SUPPLIER	ASSOCIATE (Employee)
- First name	- Company name	- First name
- Last name	- Contact first name	- Last name
- Address line	- Contact last name	- Address line
- City	- Address line	- City
- State	- City	- State
- Zip code	- State	- Zip code
- Phone number	- Zip code	- Phone number
- Fax number	- Phone number	- Pager number
- Tax id	- Cell Phone Number	- Social Security #
- First order date	- Fax number	- Email address
- DUNS #	- Credit Rating	- Hire date
	- First PO Date	- Clock #
	- DUNS #	

An ASSOCIATE can be assigned to several CUSTOMERs,
and manage the relationship with many SUPPLIERs.
A CUSTOMER or SUPPLIER can contact multiple ASSOCIATES.

N



Attribute Generalization - Financial

SSTYPE

12 Suppose you saw a table defined like this:

How many rows would it have?

What would YOU want to do?

FINANCIAL DATA:

Dept	Year	Qtr	Bud/Act	Category	Amount
} Identifier					

A Classic Fact Table for a Dimensional Model!

How many rows would this table have?

FINANCIAL DATA:

*Dept
*Year

Qtr1 Budget Material Amount
Qtr2 Budget Material Amount
Qtr3 Budget Material Amount
Qtr4 Budget Material Amount
Qtr1 Budget Labor Amount
Qtr2 Budget Labor Amount
Qtr3 Budget Labor Amount
Qtr4 Budget Labor Amount
Qtr1 Budget Capital Amount
Qtr2 Budget Capital Amount
Qtr3 Budget Capital Amount
Qtr4 Budget Capital Amount
Qtr1 Actual Material Amount
Qtr2 Actual Material Amount
Qtr3 Actual Material Amount
Qtr4 Actual Material Amount
Qtr1 Actual Labor Amount
Qtr2 Actual Labor Amount
Qtr3 Actual Labor Amount
Qtr4 Actual Labor Amount
Qtr1 Actual Capital Amount
Qtr2 Actual Capital Amount
Qtr3 Actual Capital Amount
Qtr4 Actual Capital Amount



Extreme Attribute Generalization

SSTYPE
13

ENTITY
*EntityID
EntityName

ATTRIBUTE
*EntityID
*AttributeName
AttributeValue

ATTRIBUTE
*EntityID
*Name
Value
Type
Size
Precision
Units
LastUpdate
Source
Confidence
...

What is lost here?

What is hard?

What is easy?

How many rows?

What else might be of interest about an attribute?

The Power of Generalization Thinking!

If find you are mixing value domains, you may have generalized too much.

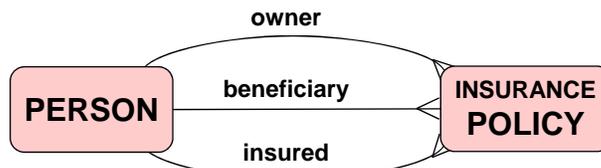
N



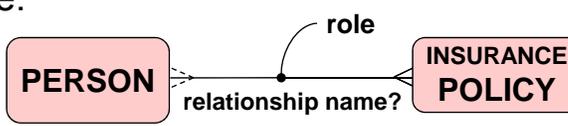

Relationship Generalization

SSTYPE
14
G. Simsion, *DM Essentials*, 2005, p.140.

Reducing multiple relationships:



to one:

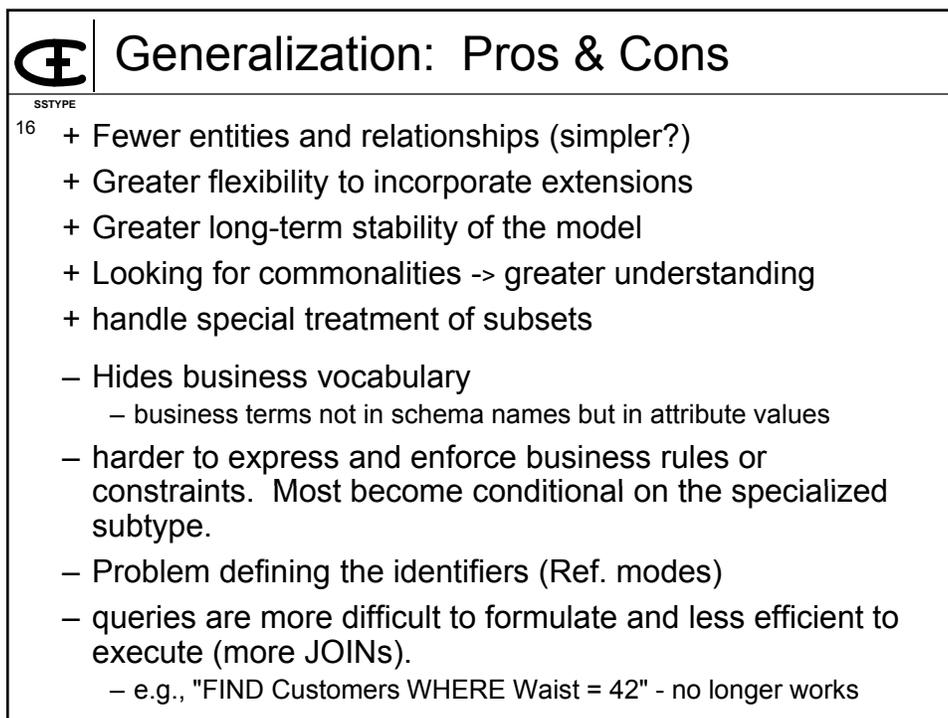
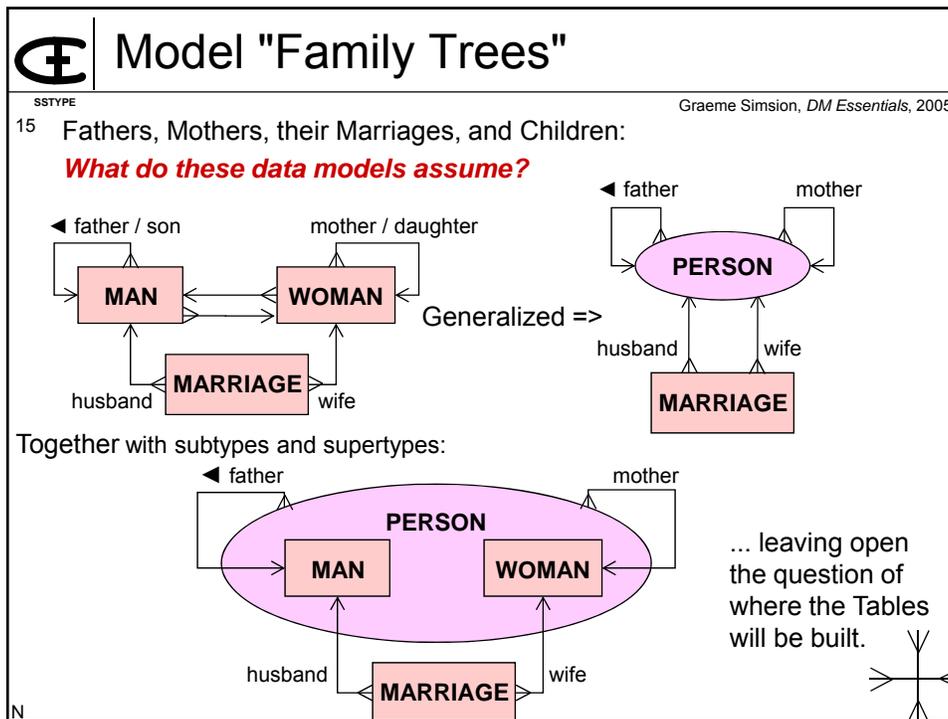


NOTE: We have already generalized the individual roles into a Person entity.

Where would you store the 'role' attribute?

How many foreign keys are stored? Where?

N



Entity Generalization



SSTYPE

17 Discerning Inter-Entity Relationships:

PERSON

ORGANIZATION

EMPLOYEE

SHAREHOLDER

CUSTOMER

VENDOR

Construct a high-level, conceptual data model.

Problems?

Observations?

N



A Fundamental Assumption in a Data Model Diagram

SSTYPE

18

- The main construct is an Entity.
- Each labeled box/circle represents an Entity Type
 - a Defined Structure (a schema template)
 - a Population of Instances
- Grouping Instances into Types is essentially Arbitrary.
 - The world isn't naturally that way; the designer imposes a view
- All Entity Type Populations are strictly Disjoint (mutually exclusive; or non-overlapping).
 - At least that is the system's assumption, thus each file/table has its own set of records/tuples.

Is this always true?

What about:

EMPLOYEE

SHAREHOLDER

E

Using Subtypes and Supertypes

SSTYPE
Smith & Smith, ACM TODS, 1977/6.

19

- Subtypes and Supertypes allow us to formally represent overlapping populations
 - Every member of a subtype population is-a member of its supertype population(s).

so we can model:

- different roles played by members of a common population, e.g.:
 - role determines the attributes

```

            graph BT
            EMPLOYEE --> PERSON
            SHAREHOLDER --> PERSON
            CUSTOMER --> PERSON
            ORGANIZATION --- PERSON
            style ORGANIZATION stroke:#f00,stroke-width:2px
            style PERSON stroke:#f00,stroke-width:2px
            style EMPLOYEE stroke:#f00,stroke-width:2px
            style SHAREHOLDER stroke:#f00,stroke-width:2px
            style CUSTOMER stroke:#f00,stroke-width:2px
            
```

- different states of an entity (over time), e.g.:

```

            graph LR
            A[ORDER RECEIVED] -.-> B[ORDER VALIDATED & ACCEPTED]
            B -.-> C[ORDER FILLED]
            C -.-> D[BACK ORDER]
            D -.-> E[ORDER REJECTED]
            E -.-> A
            
```

E

Subtype-Supertype "Relationship"

SSTYPE

20

- Tempting to call it a "Relationship," but ...

```

            graph LR
            S[Supertype] ---|1:1| T[Subtype]
            style S stroke:#f00,stroke-width:2px
            style T stroke:#f00,stroke-width:2px
            
```

- $\text{pop}(\text{subtype}) \subseteq \text{pop}(\text{supertype})$
SUBSET of
- AND the related members in the two sets are the *same* instance
 - that is what makes it different from relationships in a traditional ER/Relational data model, where the entity type populations are (assumed) disjoint!



Generalization / Specialization

SSTYPE

21 **Forming Entity Types**

- An “arbitrary” choice made by the Database Designer
- Recognizing when to use Subtypes and Supertypes
- Think about the entity *populations* you are modeling

Two basic and distinct situations:

- **Generalization:** (bottom-up - from several to a common supertype)
 - When you observe commonalities (e.g., common attributes*) across multiple entity populations.
 - the members may actually be from the same population, the same type of ‘thing’, so define a common supertype.
- **Specialization:** (top-down - from one to subtypes)
 - When there is something special about a subset of a population
 - They have some unique attributes*
 - You want to treat them differently
 - e.g., Apply a constraint, or have some attributes mandatory

*NOTE: speaking of attributes in ORM, means roles in relationships with other objects.



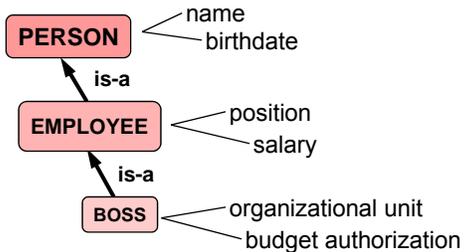
Subtypes and Supertypes



SSTYPE

22 **TWO CONDITIONS MUST ALWAYS BE TRUE:**

- each subtype population must be a **subset** (potentially) of each of its supertype populations
i.e., each instance of the subtype is in *every* supertype population
- each subtype inherits *all* the roles of its supertypes *and* must have **additional roles**/relationships



```

graph BT
    PERSON[PERSON] --- name
    PERSON --- birthdate
    EMPLOYEE[EMPLOYEE] -- is-a --> PERSON
    EMPLOYEE --- position
    EMPLOYEE --- salary
    BOSS[BOSS] -- is-a --> EMPLOYEE
    BOSS --- organizational_unit[organizational unit]
    BOSS --- budget_authorization[budget authorization]
            
```

↑

More Instances
(larger population)

↓

More Attributes
/Roles
/Relationships

If either condition is NOT true,
no reason to call out the subtype in a separate definition.

Diagramming Subtypes and Supertypes

23 Two Basic Representations:

1. **NESTED** (Euler Diagram)

- + Intuitive - visually shows inclusion
- + Clean and Compact
- Generally assumed disjoint
- Only good for simple cases
- Not good for complex cases - difficult to represent both exclusive and overlapping subtypes (like a Venn Diagram):

Diagramming Subtypes and Supertypes

24 Two Basic Representations:

2. **SEPARATED**

- + More common
- + Easier to show constraints and multiple supertypes for more complex cases.
- not visually intuitive
- confusion with "relationship"
- Adds more "clutter"

E
Diagramming Exercise
G

SSTYPE

25

- Convert the following Nested diagram into a Separated S/Type diagram:

How to model 'E' ?

Any constraints required?

What if only exclusive subtypes allowed?

E
Subtype / Supertype Constraints

SSTYPE

26

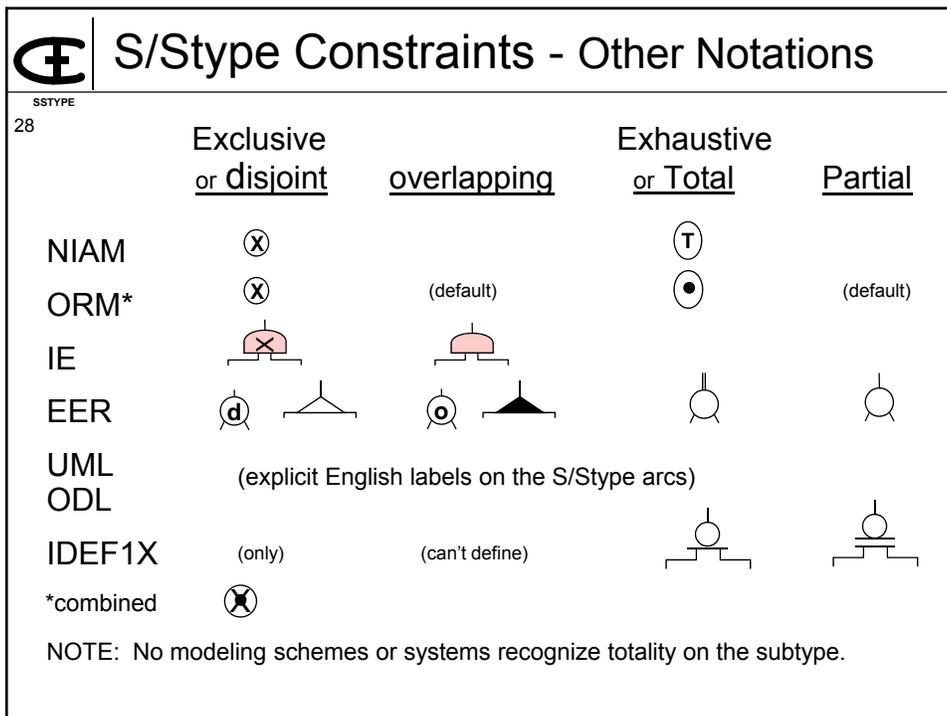
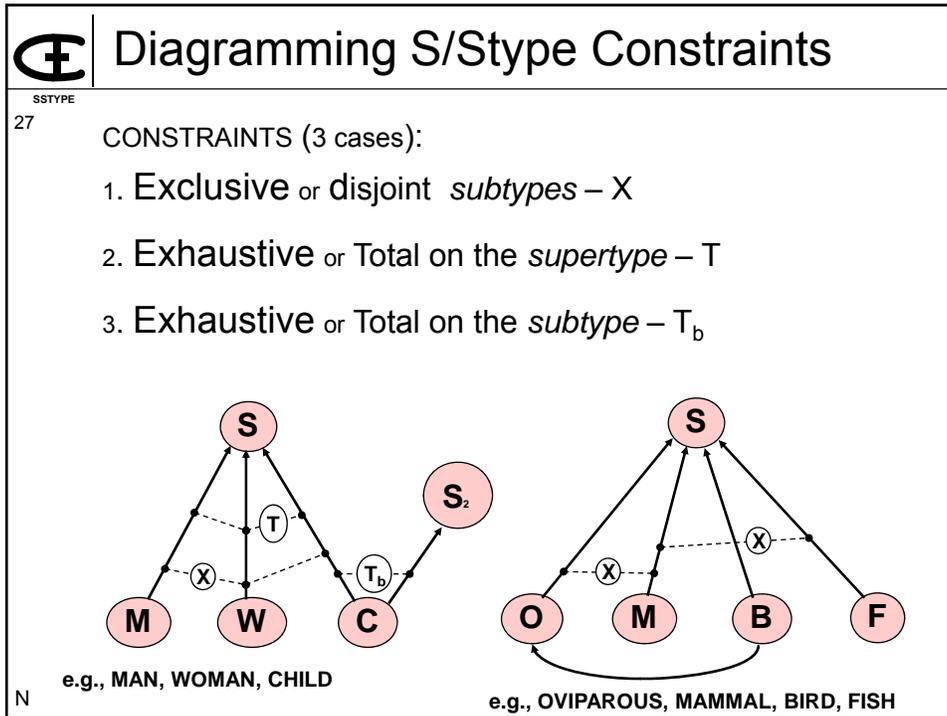
WITHOUT CONSTRAINTS,
ASSUME THE MORE GENERAL CASE:

- overlapping subtype populations
 - else Disjoint, so apply Exclusion constraint:
- non-exhaustive (Partial) on the supertype, i.e., a supertype instance need not be in *any* subtype
 - else Mandatory/Totality/Dependency constraint

=> Declare constraints on the more restrictive cases

- Some systems allow only Disjoint and Total
 - it is possible to model Overlapping, even if the system only allows Disjoint subtypes. **HOW?**
- Some systems make Disjoint and Total the defaults

N



⊆ "Well - Defined" Subtypes

SSTYPE

29

- based on an attribute of the supertype
 - called the "distinguishing" attribute
- characteristics of the relationship determine the constraints on the subtypes
 - mandatory attribute => exhaustive/totally constraint (⊆)
 - single-valued attribute => exclusive subtypes constraint (⊗)

- **What if an optional attribute?**
- **What if a multi-valued attribute (M:N relationship)?**

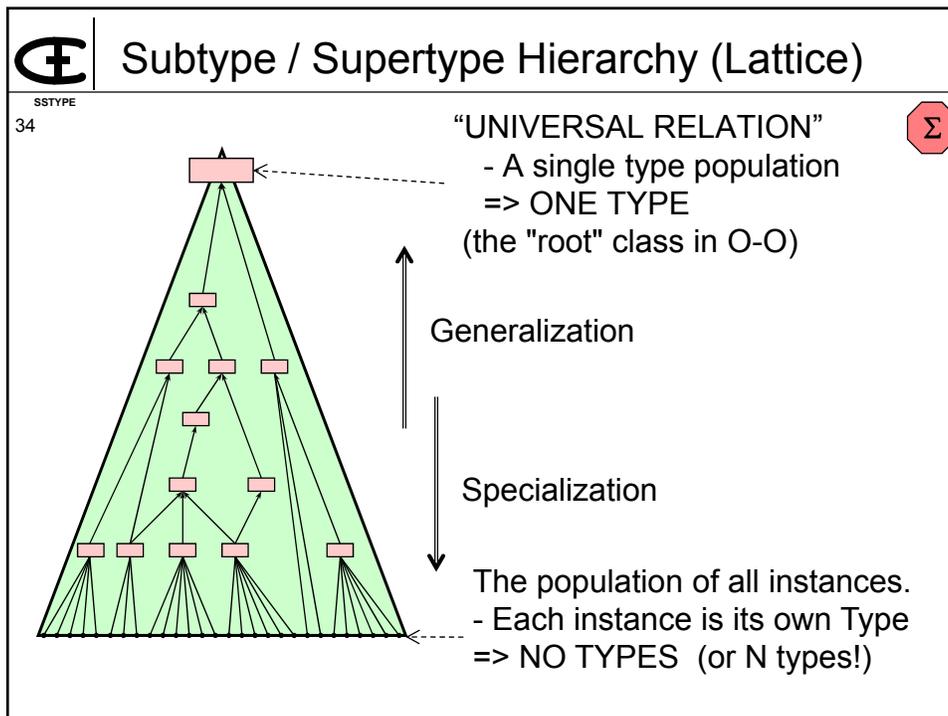
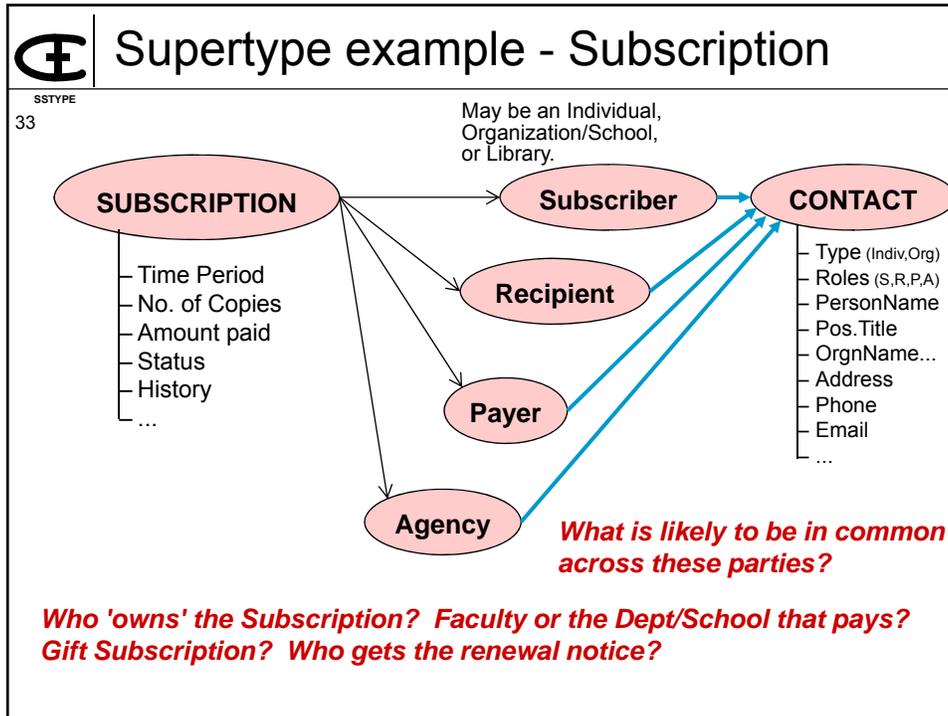
⊆ Subtype Definition

SSTYPE

30

Attribute-Defined Subtype (Intentional Set)

- a rule for including a Supertype instance in the Subtype
- Defined in terms of the values of a supertype attribute
 - in general, a Boolean expression on attribute(s) of the supertype
- can be considered a *constraint* rule on subset membership
- there are many possible subgroupings (specializations) of an entity type based on the values of its attributes, so find those that matter.
- it is not always possible to define the rules for membership in a subtype, hence: Extensional Set...





Sub/SuperType Hierarchy/Lattice - notes

SSTYPE

35 AT THE EXTREMES:

- a single supertype at the top is called the UNIVERSAL RELATION. If you built a single table for all the data in your organization:
 - what would be the entity?
 - what would be the identifier?
 - what would be the attributes?
 - would all the attributes be relevant for each row?
- at the bottom would be individual instances, each instance being its own type!
 - But sharing many attributes with other entities

The real art of database design is picking the appropriate entity types within the levels of the hierarchy.

- Allows the designer to defer choosing what tables to build

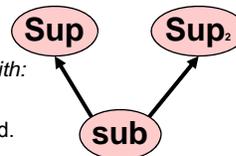


Single vs. Multiple Inheritance

SSTYPE

36

- SINGLE - every subtype has only *one* supertype
=> a strict Hierarchy of types
- MULTIPLE supertypes for a (Shared) Subtype
- If multiple supertypes, they must converge on one population higher up = the root type
=> a “partial lattice” *What's wrong/incomplete with:*



LATTICE = a partially ordered set in which every 2-element subset has both a least upper bound and a greatest lower bound.

- no lattice if no overlapping subtypes
- it is possible to transform multiple inheritance (lattice) to a generalization *hierarchy* by defining all possible combinations of subtypes
- each mini hierarchy has a root (thus partial lattice), all root objects are disjoint.

E

Extreme Entity Generalization

SSTYPE

37

© ROBIN WADE, 1990

“The level of generalization is *critical*.” - Graeme Simsion

What is the ‘THING’?

SIMSION
 & BOWLES
 ASSOCIATES

E

Comparing Modeling Schemes

SSTYPE

38

	ER/Rel (Chen/Codd)	EER (Teorey**)	ORM (Halpin)	UML (OMG)	ODL (ODMG)	SQL 1999
Class Hierarchy	X	Y	Y	Y	Y	Y
Disjoint*		Y	Y	default	only	only
Overlapping		default	default	Y	X	X
Total covering*		Y	Y	user-defined	on abstr.class	partly
Partial		Y	Y	Y (incomplete)	Y	Y
Attribute-defined discriminator		Y	'must' but...	limited (pseudo attrib)	X	X
Shared Subclasses (multiple inheritance)		Y	Y	Y	X	X

*the more restrictive case, calling for a constraint declaration.

E

Mapping to (Relational) Tables

SSTYPE
Halpin₀₈, §11.3, p.497.

39 THREE BASIC CHOICES:

- Supertype only:
(Absorption - 'flatten' up)

$\underline{K}_P D \{ P_i \} \dots \{ A_i \} \dots \{ B_i \} \dots$

- Subtypes only: (not possible in VisioEA)
(Separation - 'flatten' down)

$\underline{K}_A \{ P_i \} \dots \{ A_i \} \dots$

$\underline{K}_B \{ P_i \} \dots \{ B_i \} \dots$

- Both:
(Partition)

$\underline{K}_P D \{ P_i \} \dots$

$\underline{K}_A \{ A_i \} \dots$

$\underline{K}_B \{ B_i \} \dots$

GIVEN:

```

graph TD
    P["P  
(K)"] --> A["A"]
    P --> B["B"]
            
```

CONSIDER:

- D = Distinguishing Attribute on P
(optional)
(single-valued?)
- Exclusive: on A & B
- A & B overlapping
- Exhaustive (Total):
P in neither A or B

PROBLEMS:

- Redundancy
- Incomplete
- Querying

E

Choosing a Mapping Strategy

SSTYPE
40

MAPPING STRATEGY	Disjoint -	Overlapping	Total	- Partial	Queries
on Supertype (absorption) (flatten up)	nulls	multiple type D	+	more nulls	+
on Subtypes (separation) (flatten down)	+	redundancy	+	arbitrary P (no place for P orphans)	joins
on Both (partition)	+	+	+	+	more joins

on Sub is best if: #P >> #D

E

Inheritance and Reuse

SSTYPE

41 • Separate but related notions - often confused (Chris Date got it right)

DESIGN NOTION based on characteristics of populations:

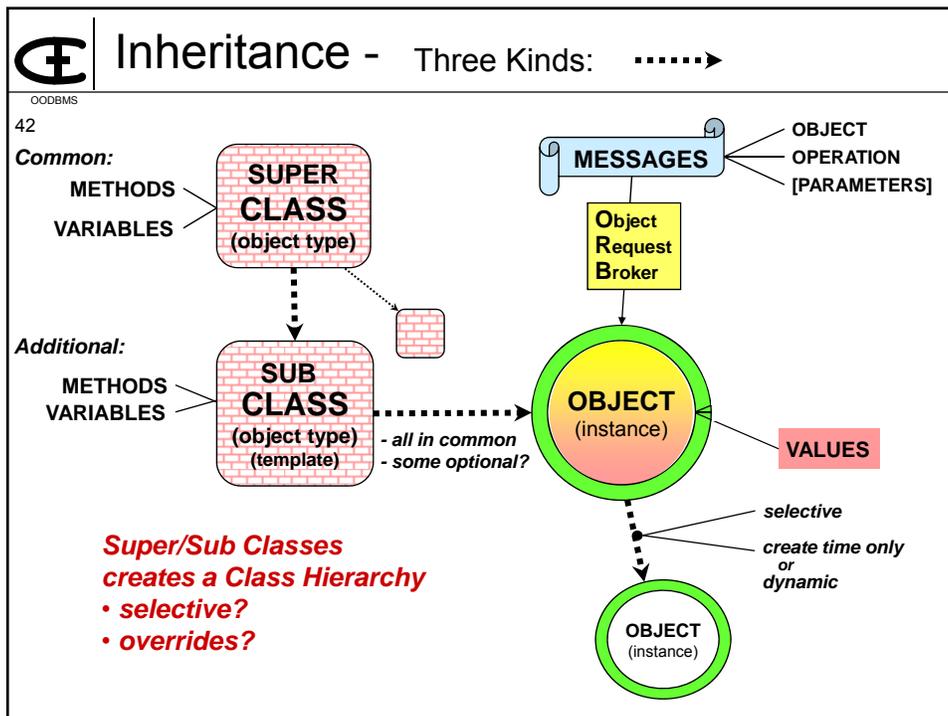
- Multiple populations with some different characteristics sharing some common characteristics, ... so define a supertype (generalization).
- Need for special treatment of a subset of a population ... so define a subtype (specialization)
 - Subtype inherits common characteristics from its supertypes plus has some additional characteristics of interest

CONSTRUCTION - implementation efficiency => REUSE

- Inheritance of definitions of data and procedures
 - for efficiency of implementation

SOLVING PROBLEMS:

- overriding and blocking
- static (copy at creation time only), vs. dynamic (maintain linkage to automatically inherit changes)
- multiple inheritance => conflict, priority order





Benefits of Subtype/Supertype

SSTYPE

G. Simsion, *DM Essentials*, 2005, p.128ff.

- 43
- Consciously and Creatively think...
 - commonalities => generalization to supertype
 - special cases => specialization to subtypes
 - Greater flexibility to handle extensions.
 - Greater stability for long-lived critical applications.
 - Generalization can reveal common patterns for reuse.
 - Abstraction for presentation, collapse the subtypes into their supertypes
 - equivalent of "leveling" in process/DFD models
 - Use subtyping to aid human understanding with no intention of implementing as separate tables.
 - Approach design top-down, bottom-up, or middle-out
 - Explicit representation of multiple table designs, thus deferring the choice for later implementation.



References

SSTYPE

- 44
- John SMITH and Diane P. SMITH, "Database Abstractions: Aggregation and Generalization," *ACM TODS*, (2:2) 1977. -classic
 - Chris J. DATE, "Subtypes and Supertypes: Setting the Scene," *Database Programming & Design*, 1999 February.
 - Terry HALPIN, *Information Modeling and Relational Databases*, Morgan Kaufmann, 2001, in ORM context.
 - Graeme C. SIMSION and Graham C. WITT, *Data Modeling Essentials*, 3e, Morgan Kaufmann, 2005, chap. 4.
 - Ramez ELMASRI and Shamkant NAVATHE, *Fundamentals of Database Systems*, 3e, Addison-Wesley, 2000, chapter 4 - academic, theoretical.
 - Steve HOBERMAN, *Data Modeler's Workbench*, Wiley, 2002, chapter 9 - practical with examples.
 - Len SILVERSTON, *The Data Model Resource Book - A Library of Universal Models for all Enterprises*, Vols. 1 & 2, Wiley, 2001.
 - Susanne W. DIETRICH and Susan D. URBAN, *An Advanced Course in Database Systems - Beyond Relational Databases*, Prentice Hall, 2005.