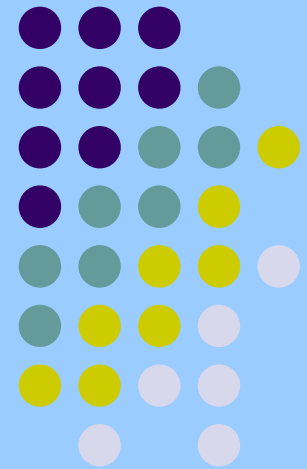
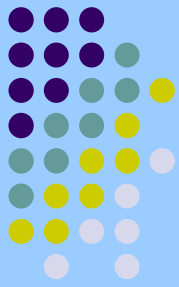


Data Modeling using XML

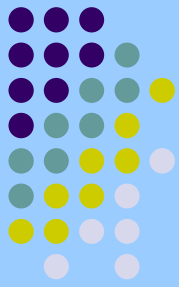
Murali Mani, WPI
Antonio Badia, University of
Louisville
Oct 13, 2003



Outline



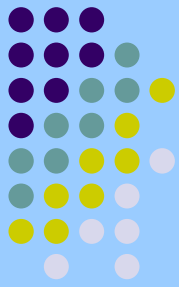
- **Part I:** How to come up with good XML designs for real world database applications?
- **Part II:** Translation between Relational and XML models.



Part I:

How to come up with good XML designs for real world database applications?

What is XML?



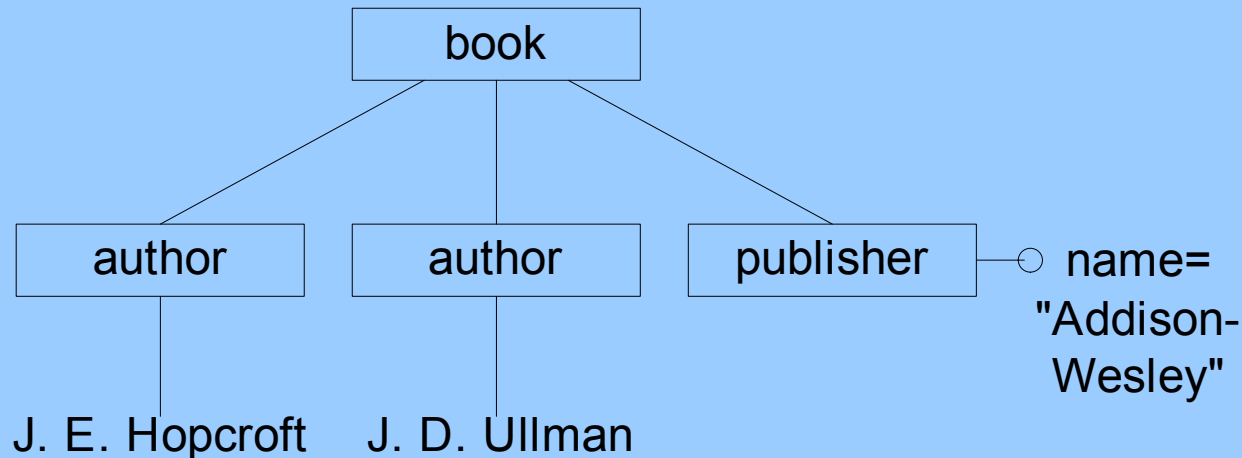
```
<book>
```

```
  <author>J. E. Hopcroft</author>
```

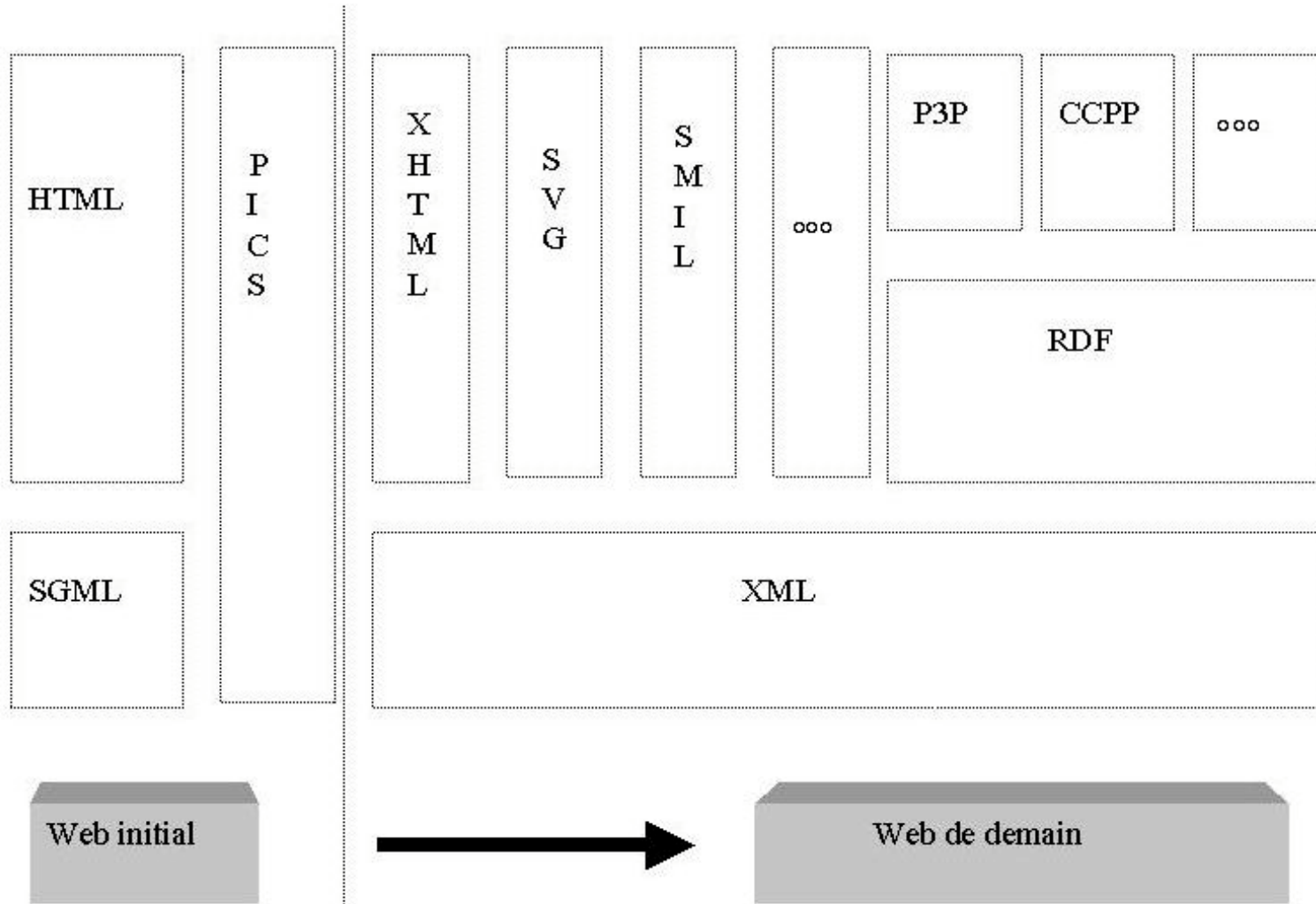
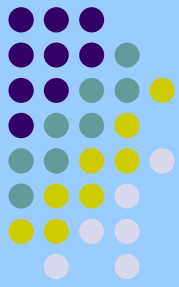
```
  <author>J. D. Ullman</author>
```

```
  <publisher name="Addison-Wesley"/>
```

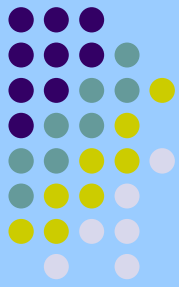
```
</book>
```



XML for information exchange



XML Publishing



Text applications

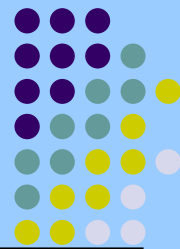
Database applications

`<reviewer>X</reviewer>`
gave `<rating>`two thumbs
up`</rating>` to `<movie>`
Fugitive, The`</movie>`

```
<person>
  <Name>A</name>
  <Age>25</Age>
  <Salary>20000</Salary>
</person>
<person>
  <Name>B</name>
  <Age>62</Age>
  <Salary>140000</Salary>
</person>
```



Publishing Relational Databases



Name	Age	Salary
A	25	20000
B	62	140000



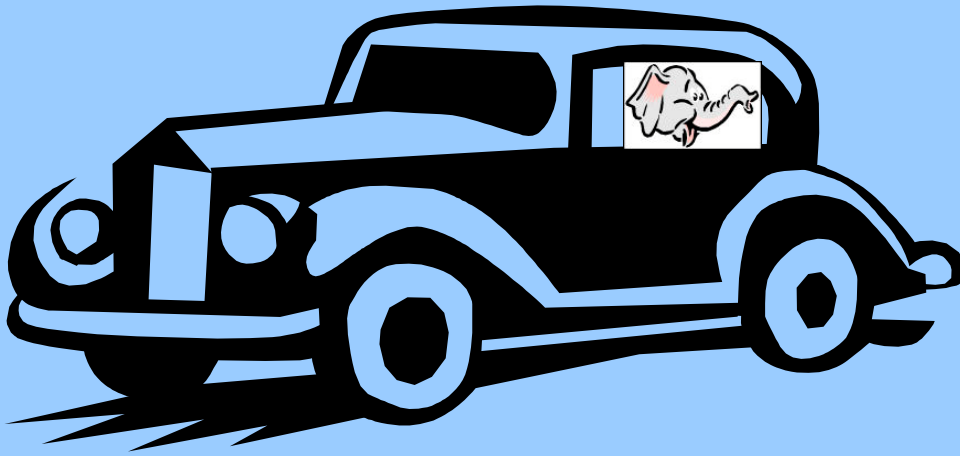
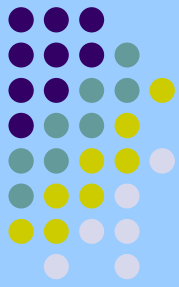
```
<person>
  <Name>A</name>
  <Age>25</Age>
  ...
</person>
```

- Users/Applications see a uniform XML view
- Exchange data with other applications
- Querying XML is easier?

Problem:

What is a good XML schema for a relational schema?

XML for Data Modeling



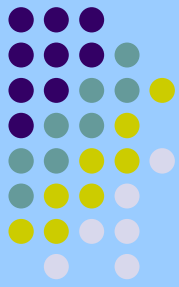
Location → location
(@val, @time, GPS)
GPS → gps (@satellite)



Location → location
(@val, @time, Bstation*)
Bstation → bstation
(@id, @sigStrength)

Location → location (@val, @time, (GPS | Bstation*))

XML as a logical data model



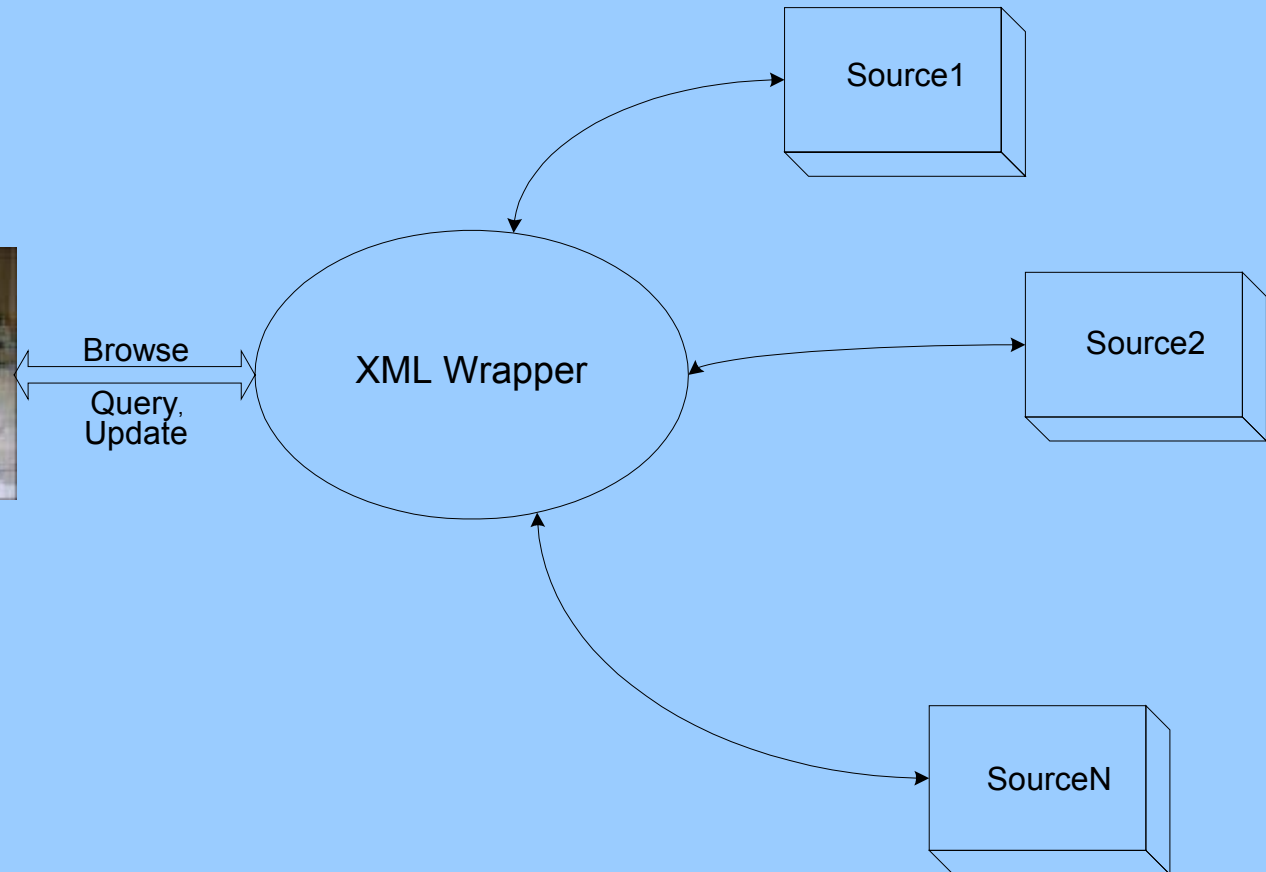
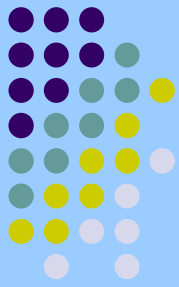
Location → location (@val, @time, (GPS | Bstation*))

- Use data modeling features provided by XML
 - Union types
 - Recursive types
 - Ordered relationships
- Easier to Query?

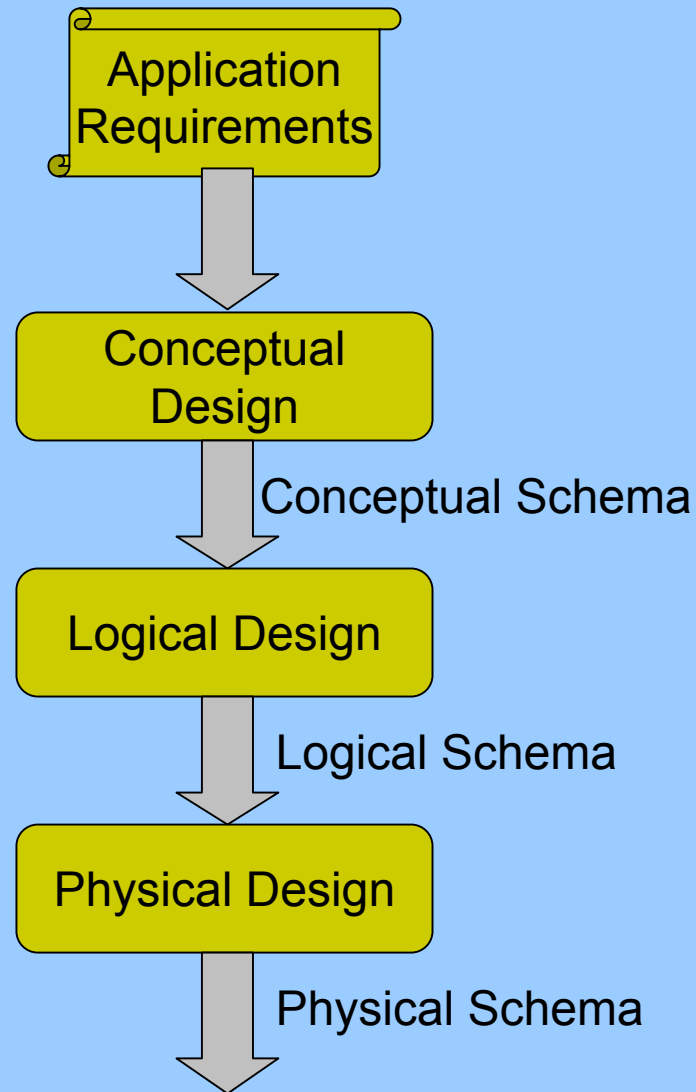
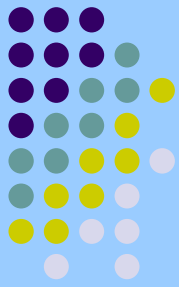
Problems:

- What is a good XML schema for an application?
- How do we store the data in relational databases?

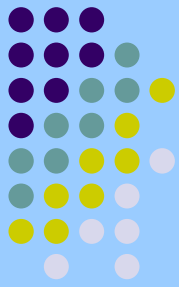
XML for data integration



Database Design Stages



Logical Data Model and Redundancy



Bad Design

Student_Professor

<u>sname</u>	BS	advisor	age
SD	CS	MXM	40
YC	EE	MXM	40

Professor

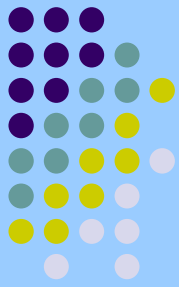
<u>pname</u>	age
MXM	40

Student

<u>sname</u>	BS	advisor
SD	CS	MXM
YC	EE	MXM

Good Design
Person

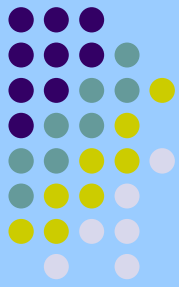
<u>name</u>	address	city	state	zip
AV	A1	Los An	CA	90034
AN	A2	Los An	CA	90034



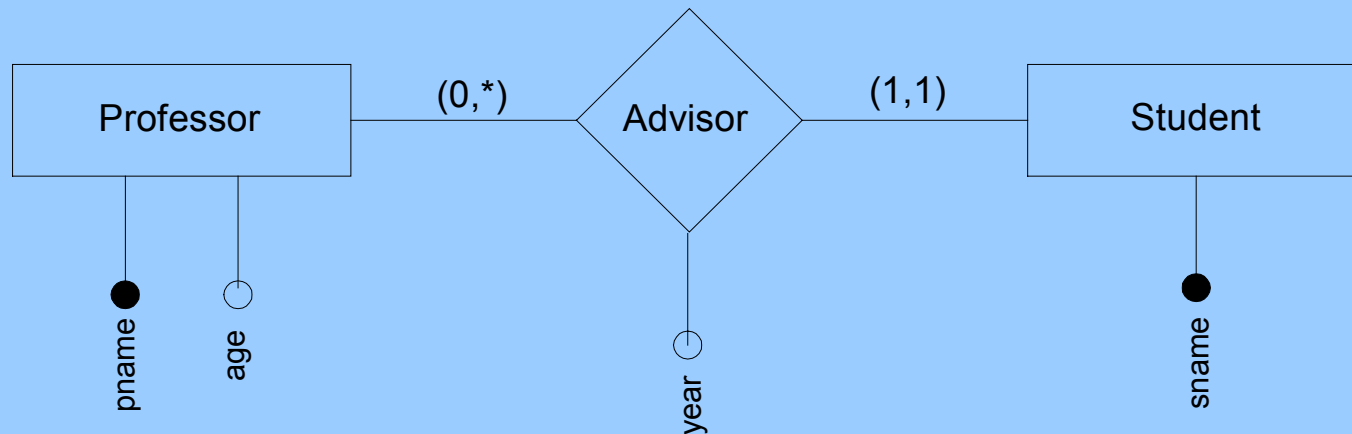
What is a Data Model?

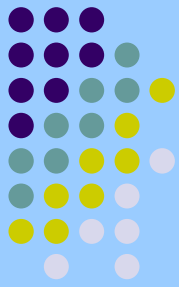
- Structural Specification
- Constraint Specification
- Operations

Entity Relationship (ER) Model



- Structures: Entity Types, Relationship Types
- Constraints: Cardinality constraints





Relational Model

- Structures: Relations
- Constraints: Key, Foreign Key

Professor	
pname	age

Student	
sname	advisor

Key Constraints:

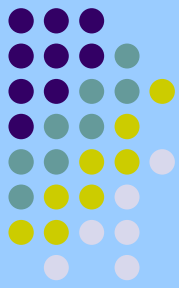
Key (Professor) = <pname>

Key (Student) = <sname>

Foreign Key Constraints:

Student (advisor) references
Professor (pname)

Specifying Structures for XML



$G = (N, T, P, S)$

$N = \{\text{Book, Author, Publisher, \#PCDATA}\}$

$T = \{\text{book, author, publisher, pcddata}\}$

$S = \{\text{Book}\}$

$\text{Book} \rightarrow \text{book} (\text{Author} +, \text{Publisher})$

$\text{Author} \rightarrow \text{author} (\text{\#PCDATA})$

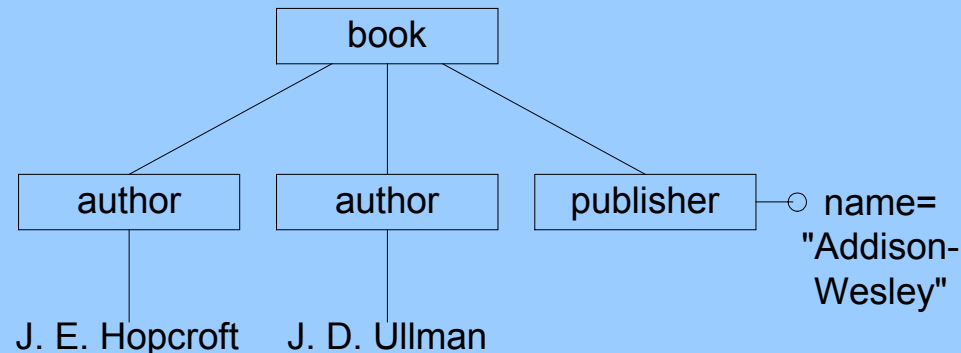
$\text{Publisher} \rightarrow \text{publisher} (@\text{name}::\text{String})$

$\text{\#PCDATA} \rightarrow \text{pcdata} (\epsilon)$

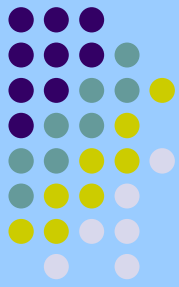
Regular Tree Grammar

Every production rule is of the form $A \rightarrow a X$

$A \in N, a \in T, X$ is a regular expression over N

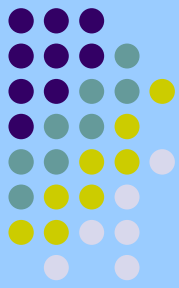


XML Schema Language Proposals



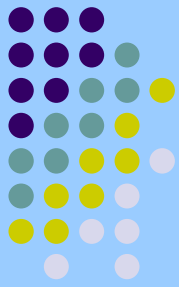
- W3C DTD: local tree grammar
- W3C XML Schema: single type tree grammar
- ISO/OASIS RELAX NG: full-fledged regular tree grammar

Properties of different Regular Tree Grammar classes



- Expressiveness
 - Regular tree grammar ***strictly more expressive*** than single type tree grammar
 - Single type tree grammar ***strictly more expressive*** than local tree grammar
- Closure properties
 - Regular tree grammar closed under union, intersection and difference
 - Single type tree grammar/local tree grammar closed only under intersection
- Type assignment
 - Type assignment can be ambiguous for regular tree grammar.
 - Type assignment is unambiguous for local tree grammar/single type tree grammar.

Ambiguous Type Assignment



$G = (N, T, P, S)$

$N = \{\text{Book, Author1, Author2, Publisher, \#PCDATA}\}$

$T = \{\text{book, author, publisher, pcddata}\}$

$S = \{\text{Book}\}$

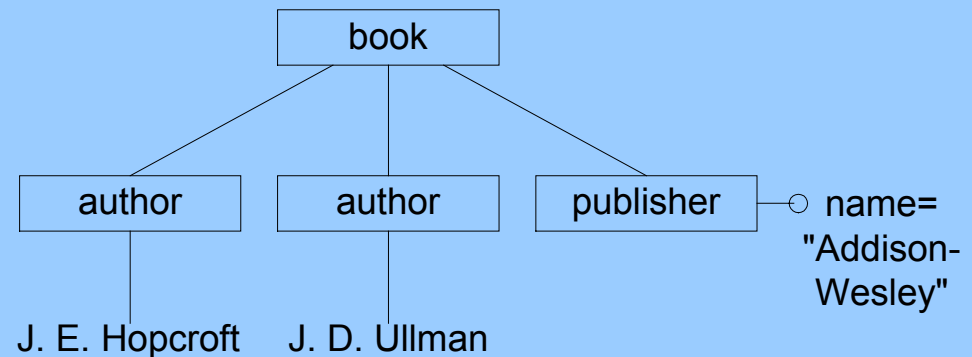
Book \rightarrow book (Author1*, Author2*, Publisher)

Author1 \rightarrow author (#PCDATA)

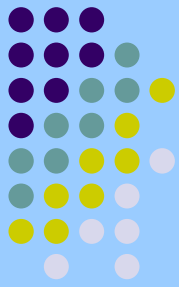
Author2 \rightarrow author (#PCDATA)

Publisher \rightarrow publisher (@name::String)

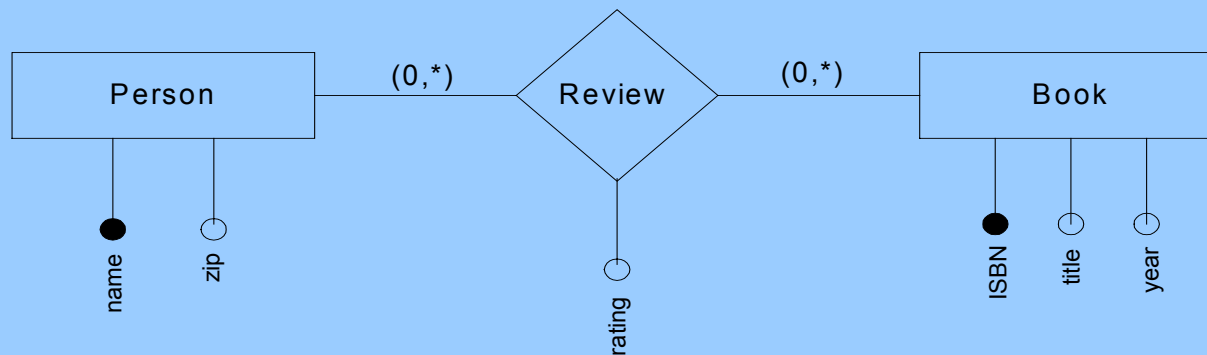
#PCDATA \rightarrow pcddata (ϵ)



Constraint Specification for XML – why?



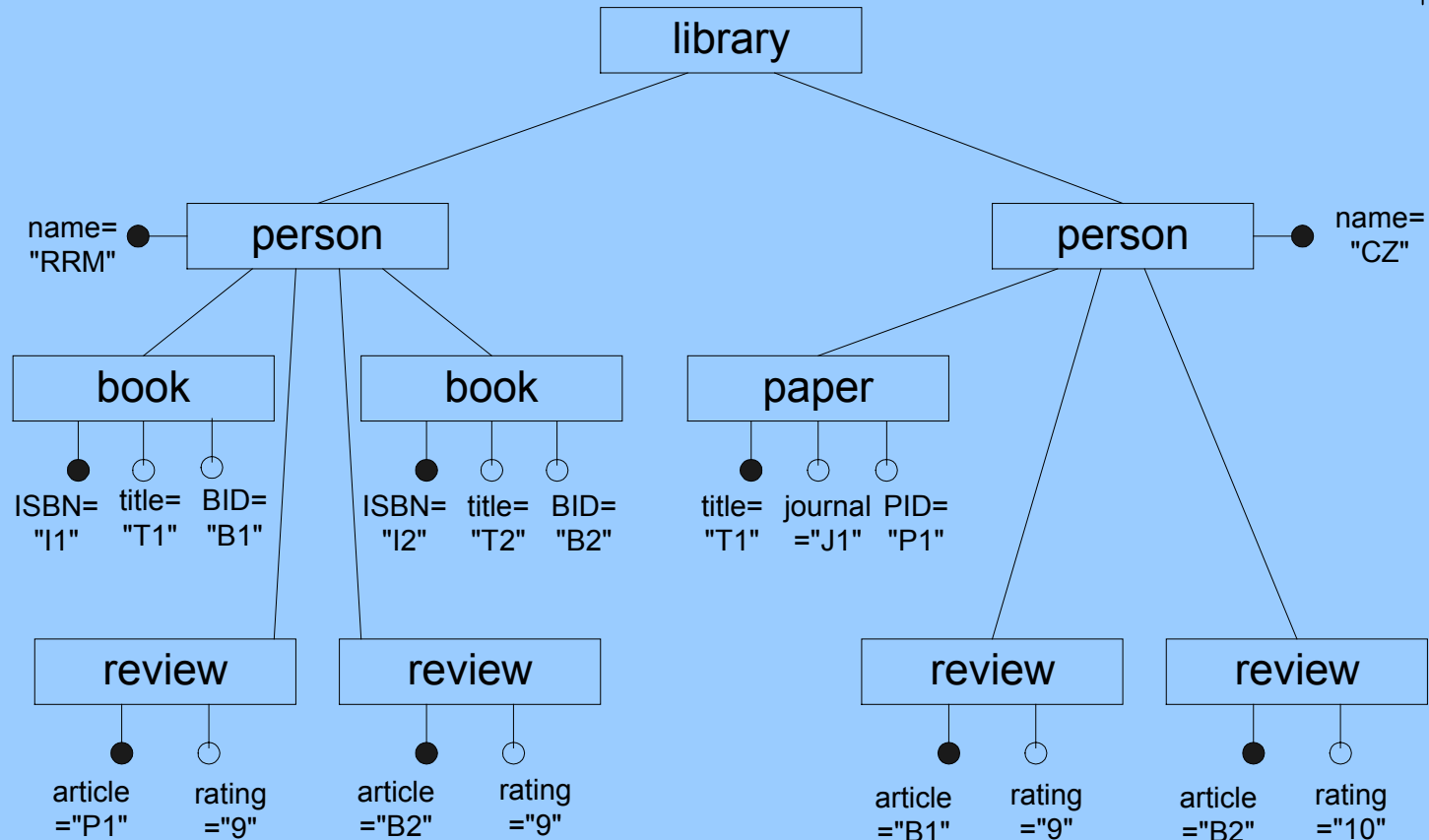
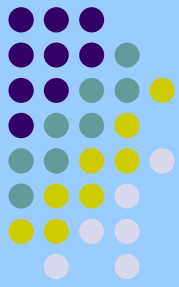
If we represent all relationships only by hierarchies, then the logical model will have redundancy.



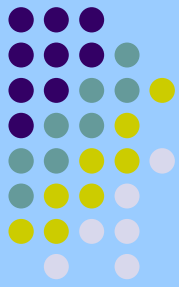
What constraint specification?

- Key, Foreign Key
- ID/IDREF

Specifying Constraints for XML: Example

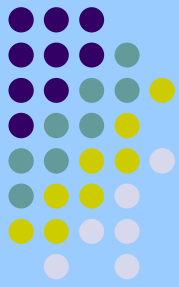


Specifying Constraints for XML

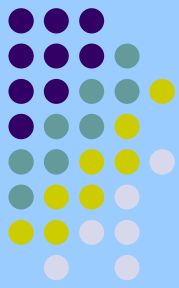


- Keys are specified using (rel, sel, field)
 - rel is relative axis
 - sel is selector axis
 - field is a set of path expressions
- For any element that “belongs to” rel, “sel” will give a set of elements. For this set of elements, field is the key.
- rel and sel can be types or path expressions
- Foreign keys are specified as (rel₁, sel₁, field₁) references (rel₂, sel₂, field₂)

Constraint Specification Proposals



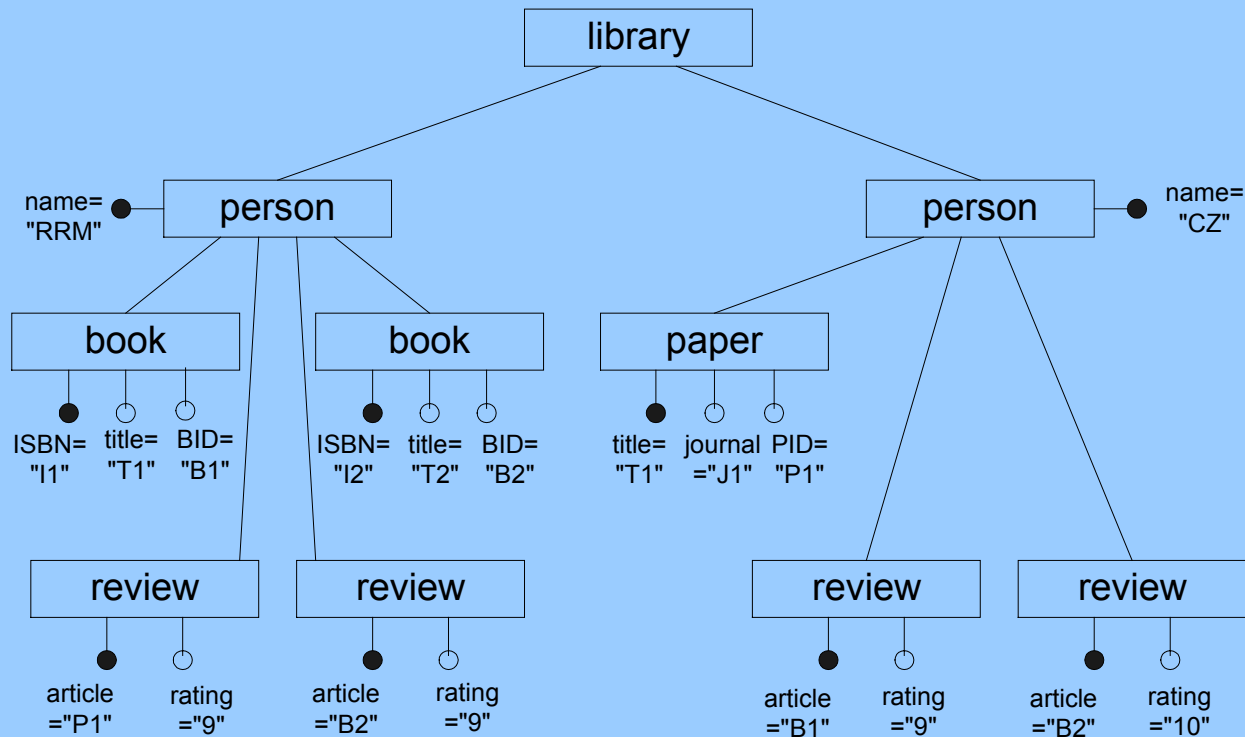
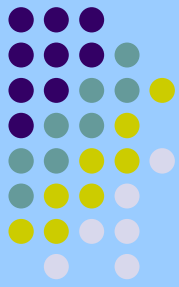
- W3C XML Schema
 - Relative axis = type
 - Selector axis = path expression
- Keys for XML – WWW10
 - Relative axis = path expression
 - Selector axis = path expression
- UCM – WWW10
 - No relative axis
 - Selector axis = type



Our proposal

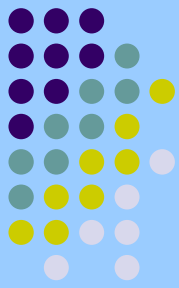
- **Relative axis = type**
- **Selector axis = type**
- **IDREF and IDREFS identify target types**

Example



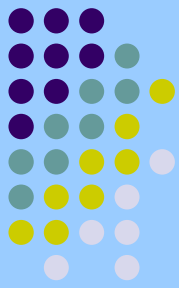
(Library, Person, <@name>
(Library, Book, <@ISBN>
(Library, Paper, <@title>
(Person, Review, <@article>)

@article::IDREF references
(Book | Paper)



Why use XML as logical data model?

Path Expressions vs Joins

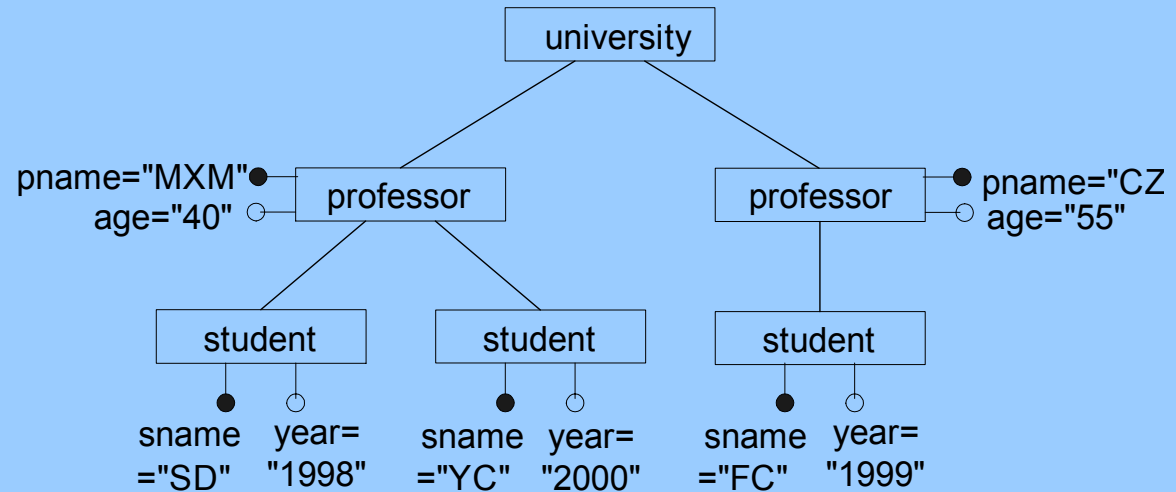


Professor

<u>pname</u>	age
MXM	40
CZ	55

Student

<u>sname</u>	year	advisor
SD	1998	MXM
YC	2000	MXM
FC	1999	CZ



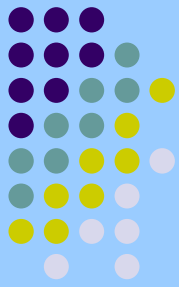
Student (advisor) references
Professor (pname)

Query: Give names of students of professors of age 40

$\pi_{\text{name}} ((\sigma_{\text{age}=40} (\text{Professor})) \otimes \text{Student})$

professor [@age=40]/student/@sname

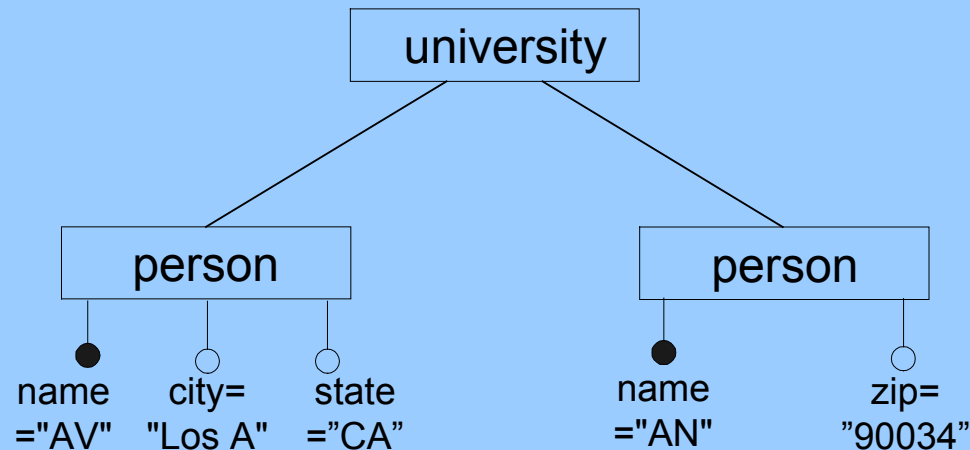
Union Types - attributes



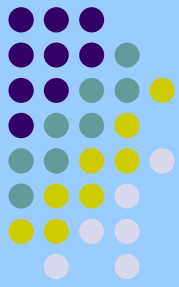
Person

<u>name</u>	city	state	zip
AV	Los A	CA	null
AN	null	null	90034

Person \rightarrow person (@name, ((@city, @state) | @zip))



Union Types - Relationships



Person

name	zip
RRM	01609
CZ	90095

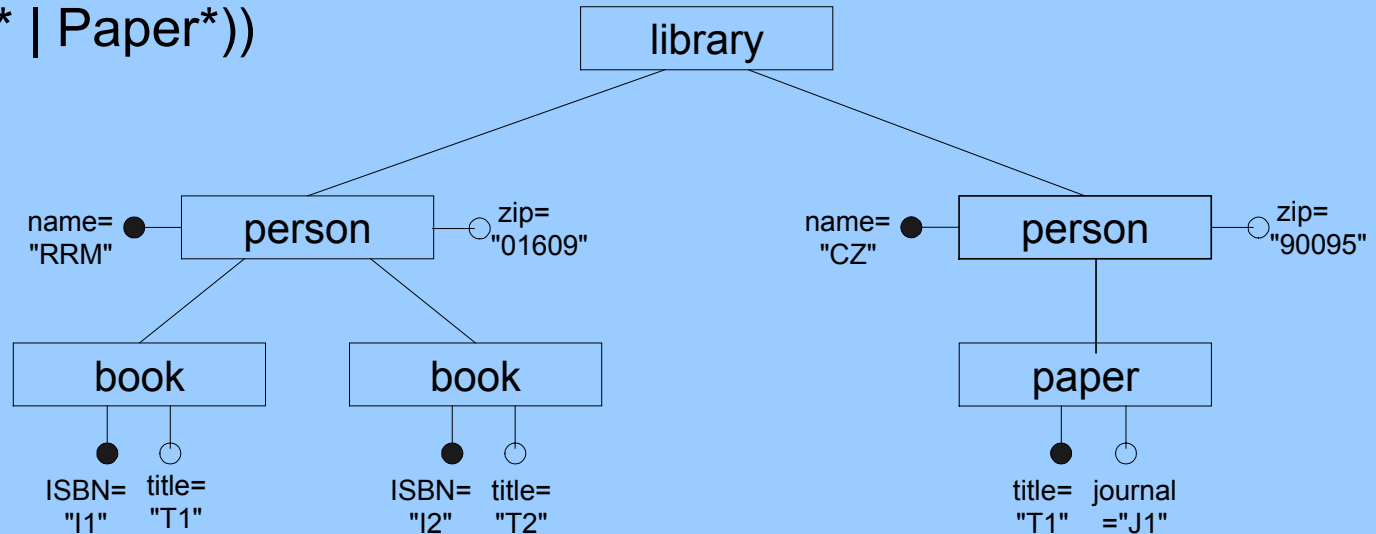
Book

ISBN	title	author
I1	T1	RRM
I2	T2	RRM

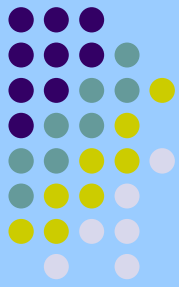
Paper

title	journal	author
T1	J1	CZ

Person → person (@name, @zip, (Book* | Paper*))



Union Types - Relationships



Conference

name	venue
ER	Chicago

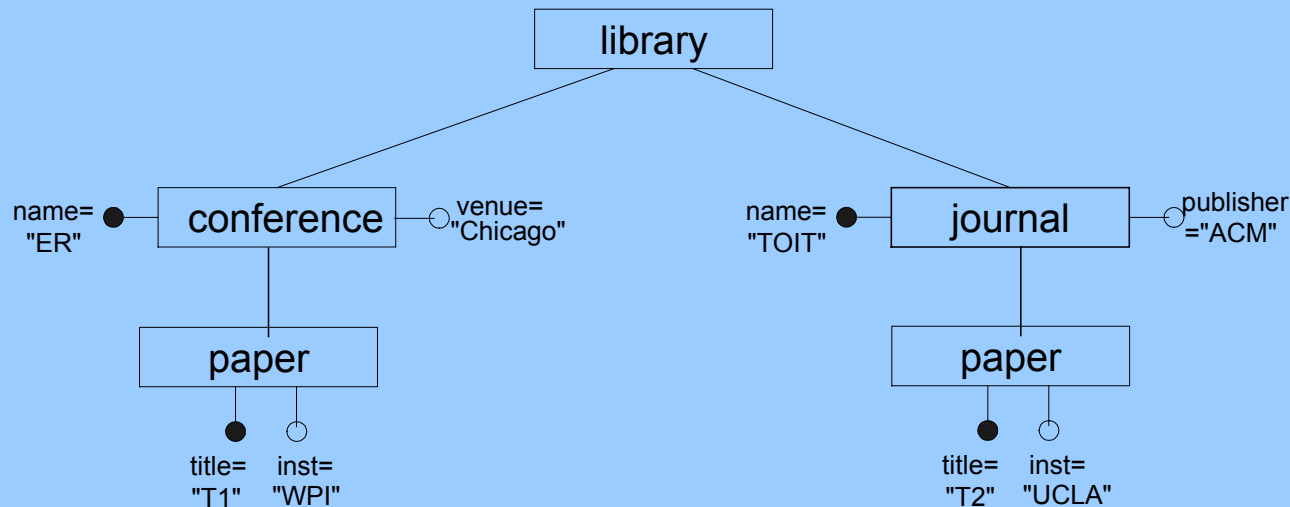
Paper

title	inst	conf	journal
T1	WPI	ER	null
T2	UCLA	null	TOIT

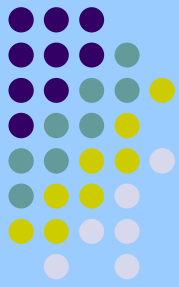
Journal

name	publisher
TOIT	ACM

Conference → conference (@name, @venue, Paper*)
 Journal → journal (@name, @publisher, Paper*)

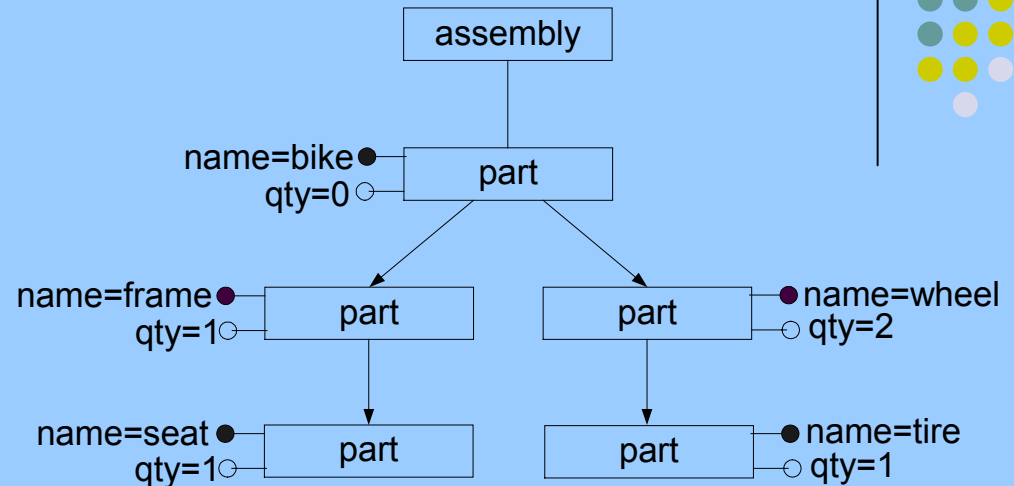


Recursive Types



Assembly

<u>name</u>	superPart	qty
seat	frame	1
tire	wheel	1
frame	bike	1
wheel	bike	2
bike	null	0

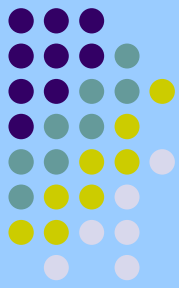


Query: What are subparts of bike?

```
WITH RECURSIVE SubPart (name) AS
  (SELECT name FROM Assembly
   WHERE superPart=bike)
UNION
  (SELECT R2.name
   FROM SubPart R1, Assembly R2
   WHERE R2.superPart = R1.name)
SELECT * FROM SubPart
```

part[@name=bike]//part/@name

IDREF vs Foreign Keys

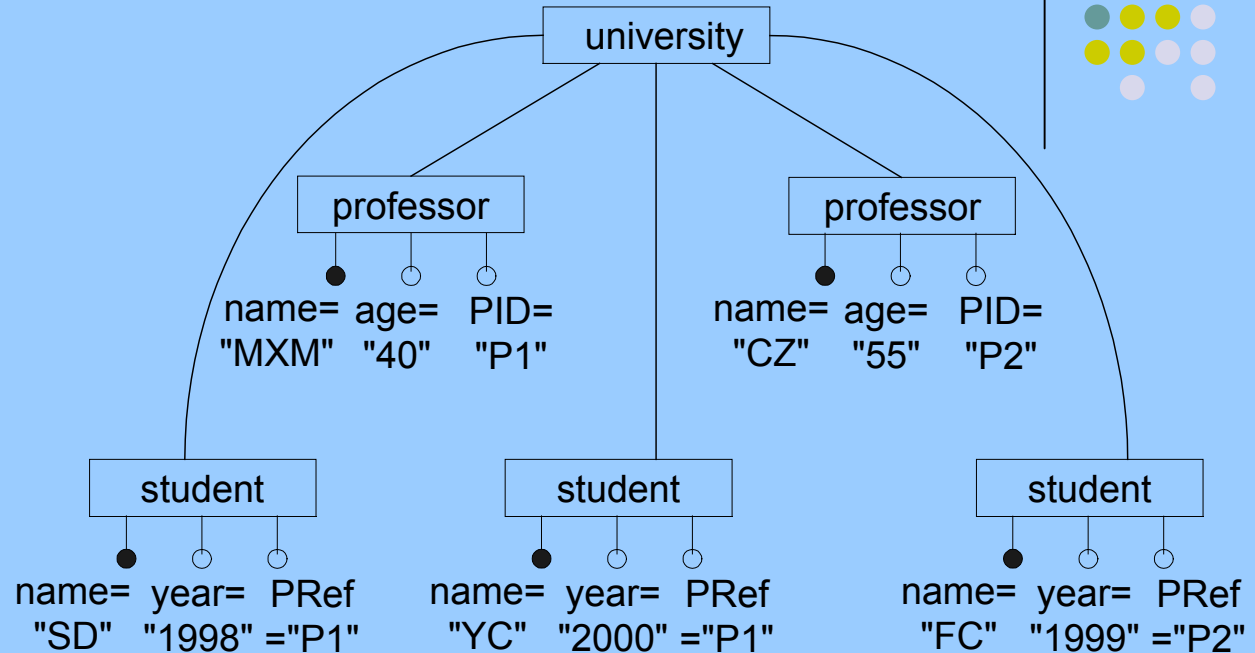


Professor

<u>name</u>	age
MXM	40
CZ	55

Student

<u>name</u>	year	advisor
SD	1998	MXM
YC	2000	MXM
FC	1999	CZ



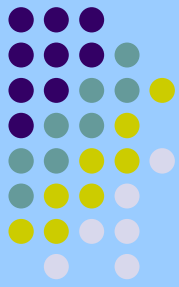
Student (advisor) references
Professor (name)

@PRef::IDREF references (Professor)

Query: Give names of students of professors of age 40

`student[@PRef=>professor/@age=40]/@name`

IDREF as union of foreign keys



Book

ISBN	title
I1	T1
I2	T2

Person

name	zip
RRM	01609
CZ	90095

Paper

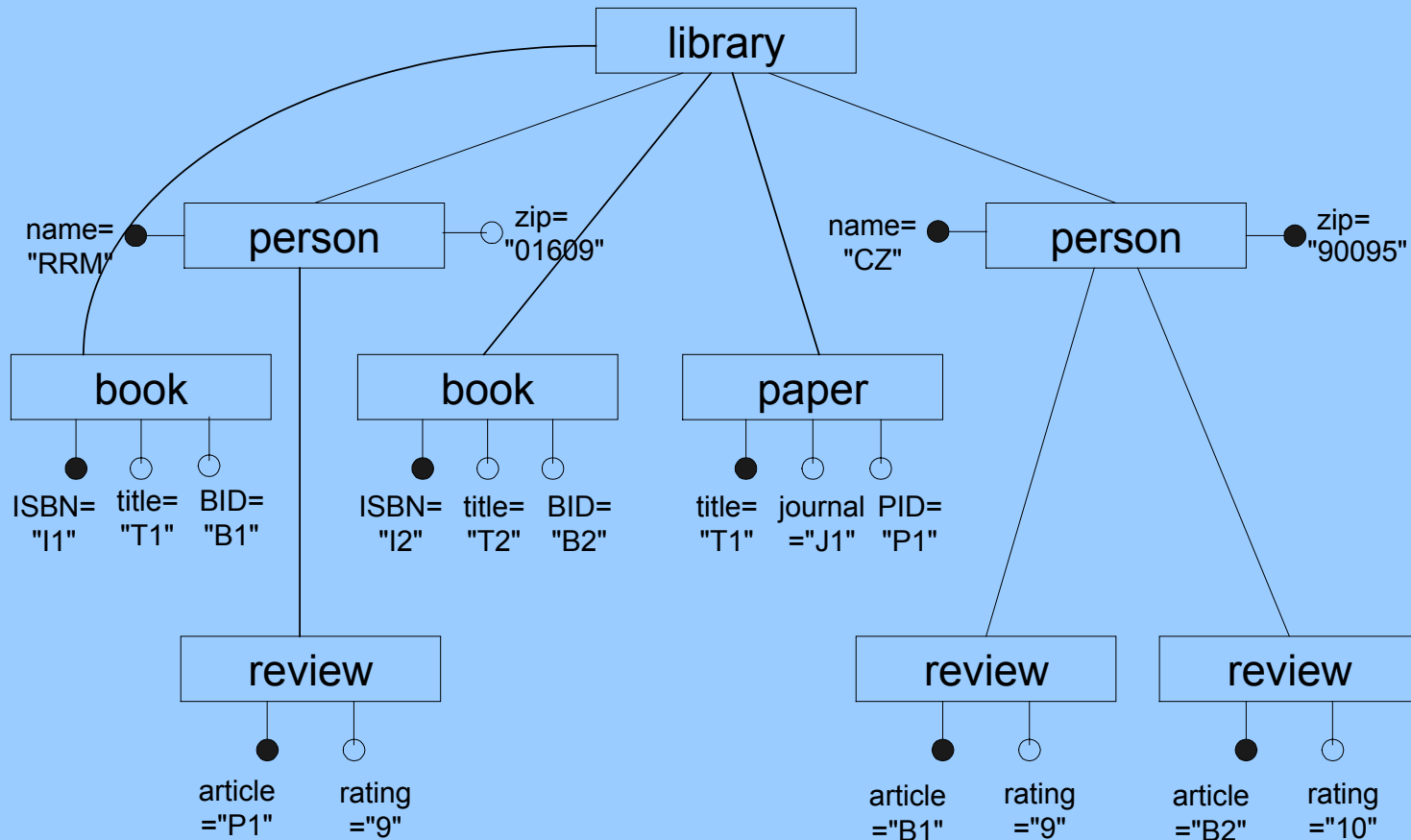
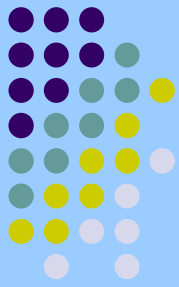
title	journal
T1	J1

Review

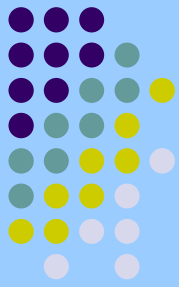
name	book	paper	rating
RRM	null	T1	9
CZ	I1	null	9
CZ	I2	null	10



IDREF as union of foreign keys



@article::IDREF references (Book | Paper)



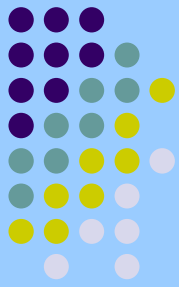
Conceptual Model: **ERex** (ER extended for XML)

Oct 13, 2003

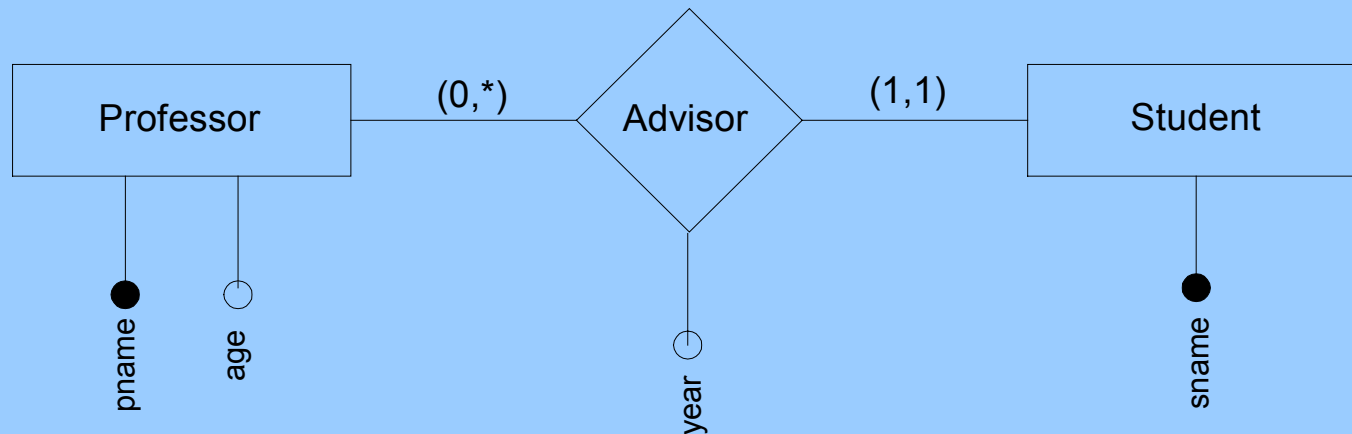
Murali Mani, Antonio Badia



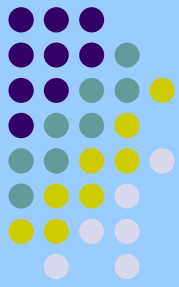
Entity Relationship (ER) model



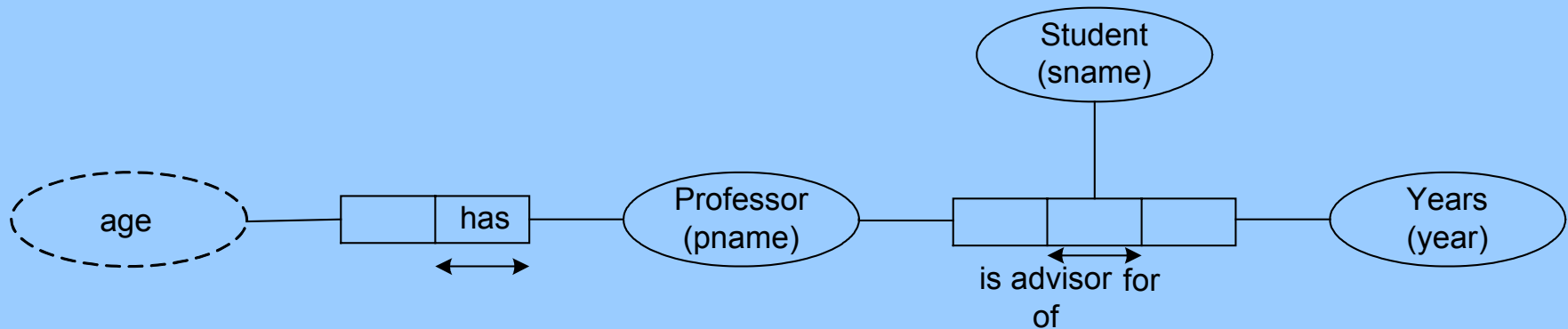
- Entity Types, Relationship Types and their attributes



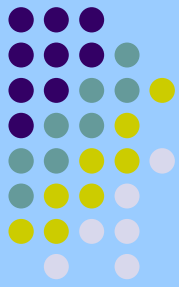
Object Role Modeling (ORM)



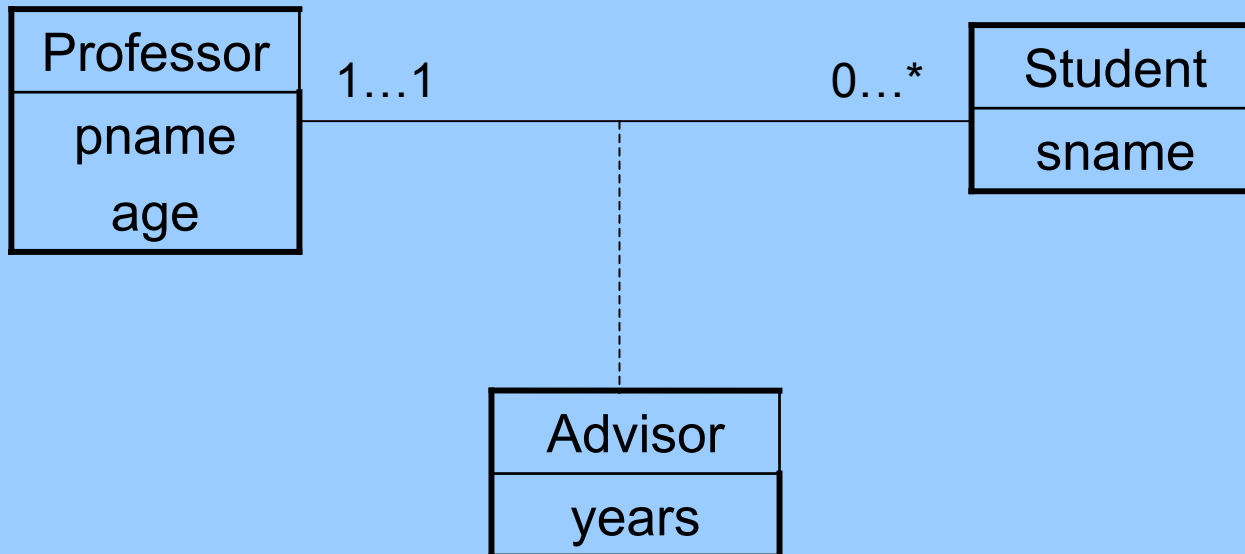
- closer to natural language sentences
- attributes/relationships are expressed uniformly using roles



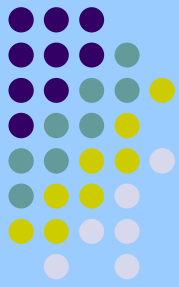
Unified Modeling Language (UML)



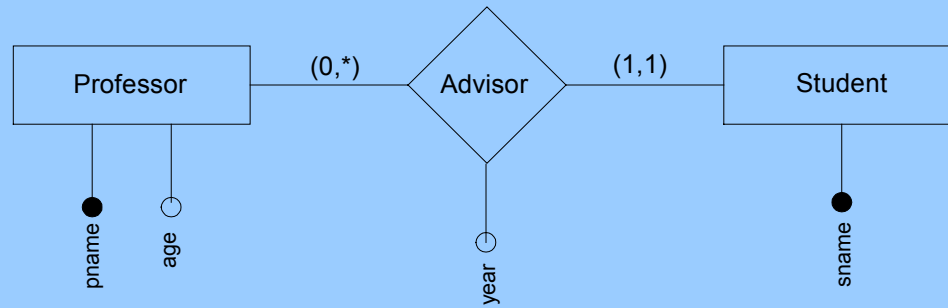
- Modeling software systems
- Class Diagrams, Association Classes



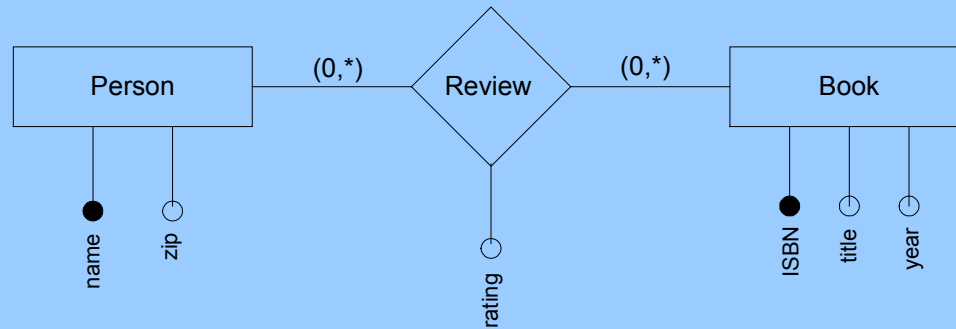
From ER model



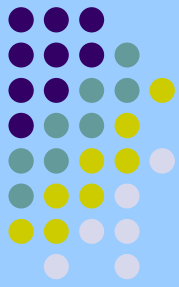
Binary 1:n relationships



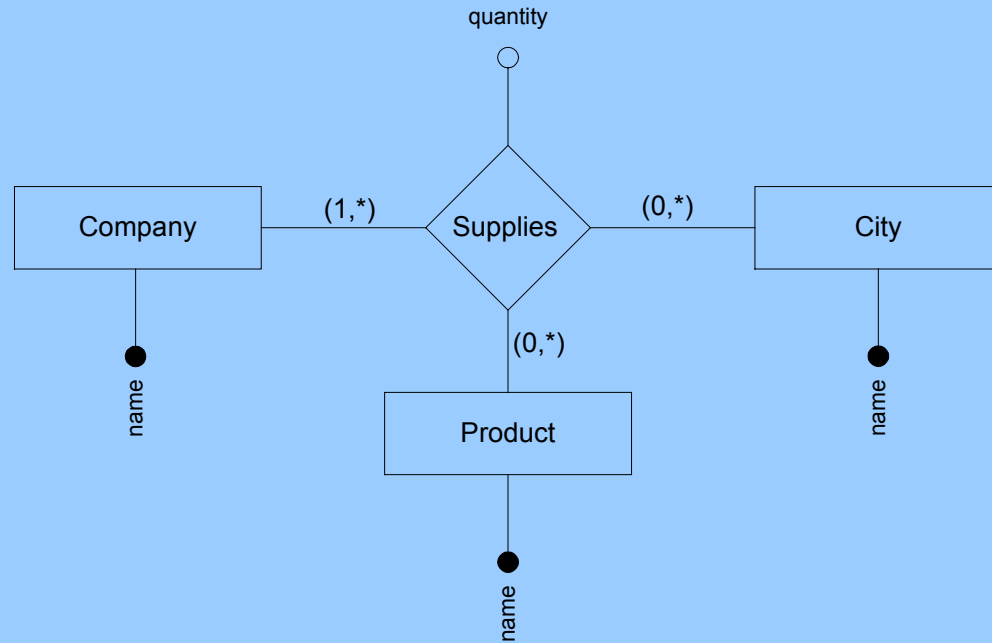
Binary m:n relationships



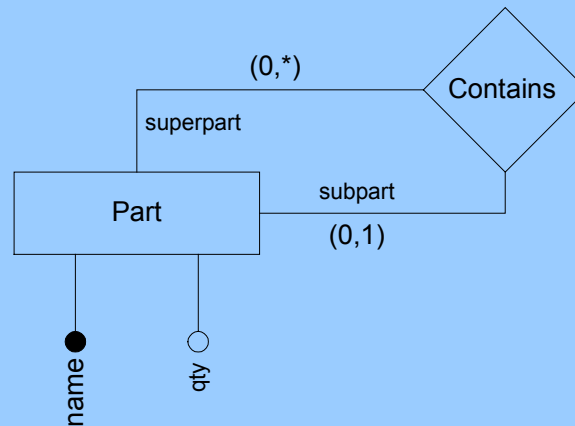
From ER Model



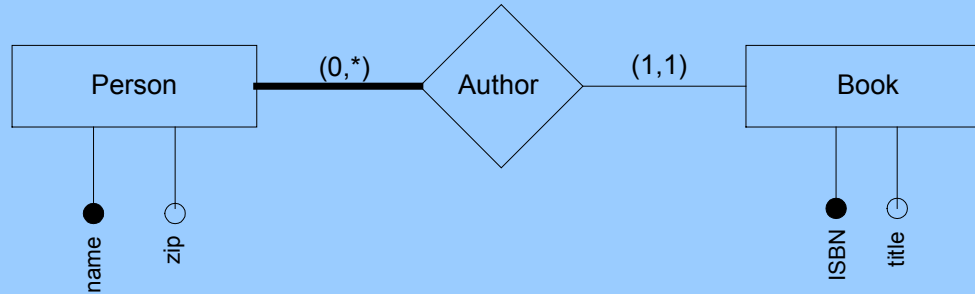
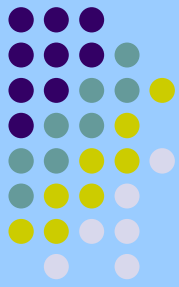
N-ary relationships



Recursive relationships



Ordered Relationships



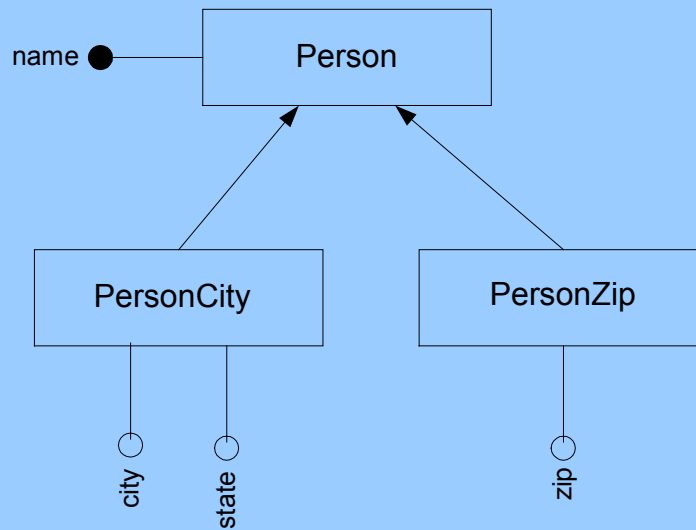
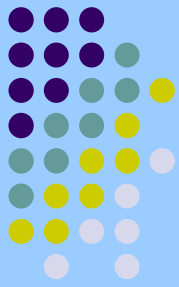
Person

<u>name</u>	zip
Ullman	95123
RRM	90095

Book

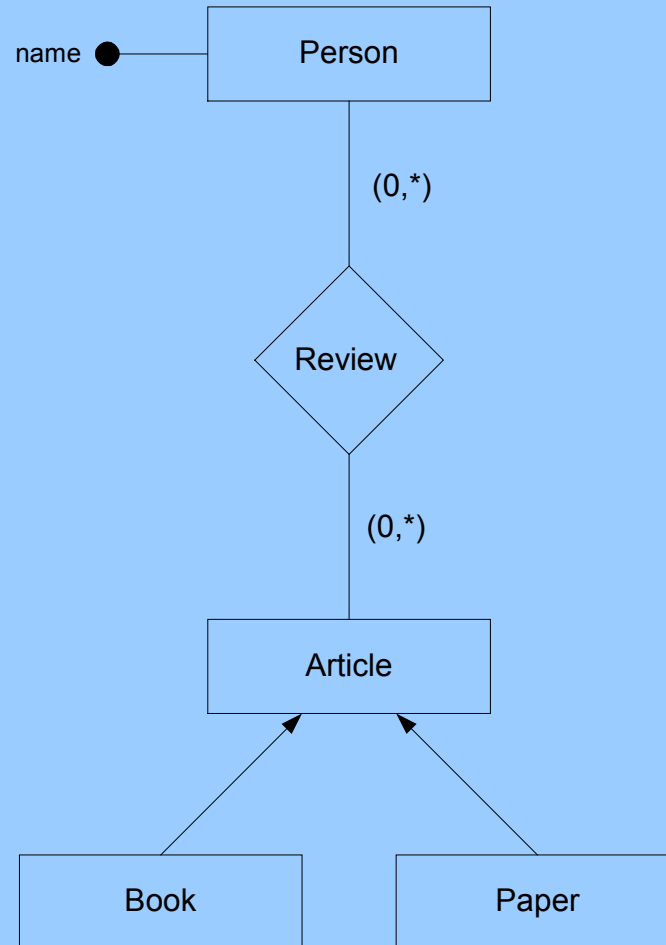
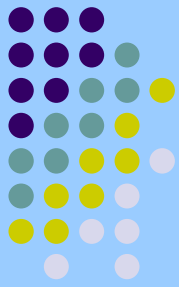
<u>ISBN</u>	title	order	author
B1	DB	2	Ullman
B2	Aut	1	Ullman
B3	MMSL	1	RRM

Categories and Set Constraints

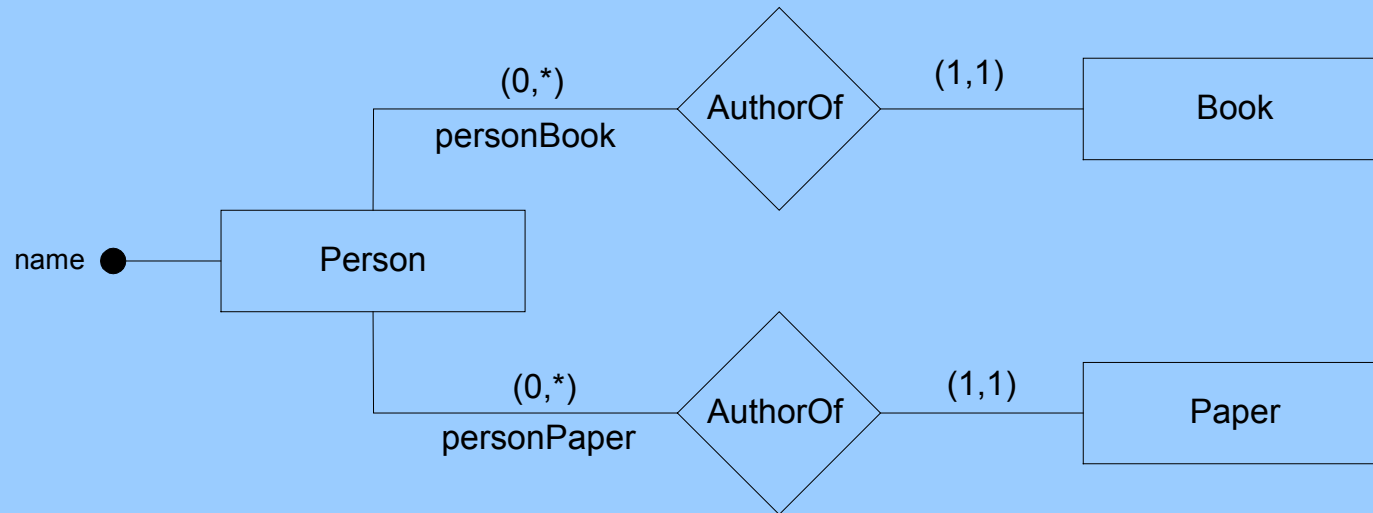
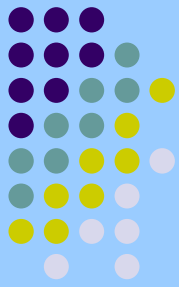


$PersonCity \cap PersonZip = \emptyset$
 $PersonCity \cup PersonZip = Person$

Categories

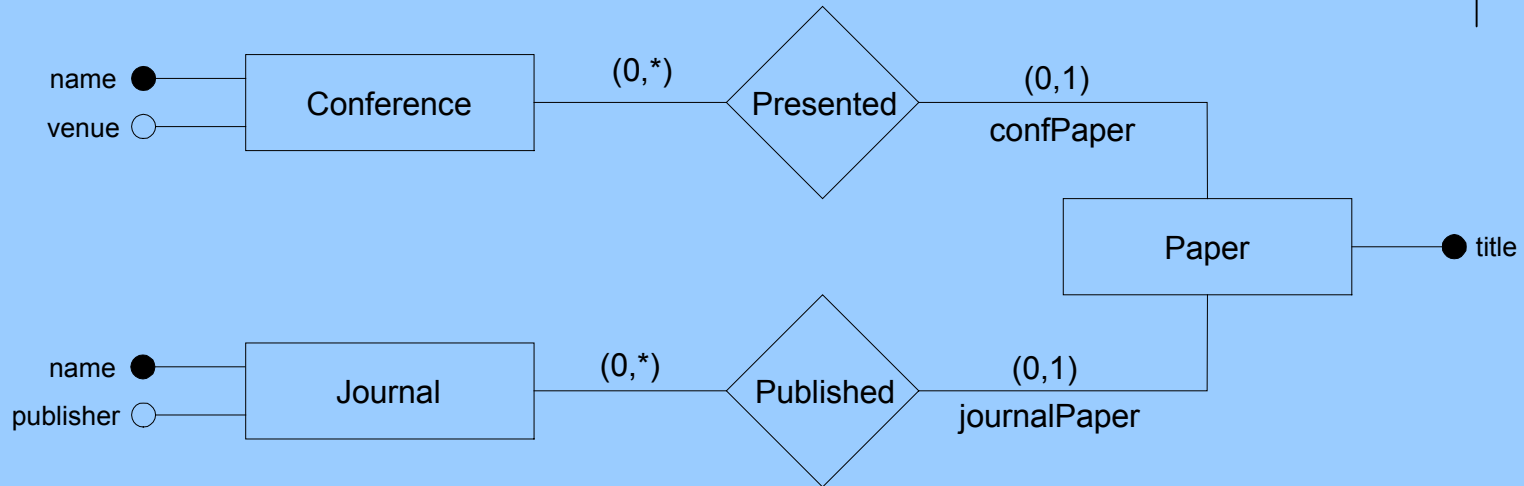
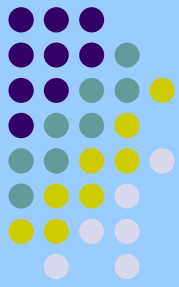


Set constraints on Roles

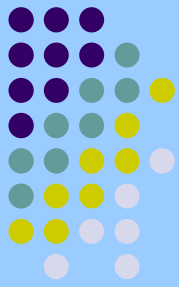


$$\text{personBook} \cap \text{personPaper} = \emptyset$$

Set Constraints on Roles



$$\text{confPaper} \cap \text{journalPaper} = \emptyset$$
$$\text{confPaper} \cup \text{journalPaper} = \text{Paper}$$



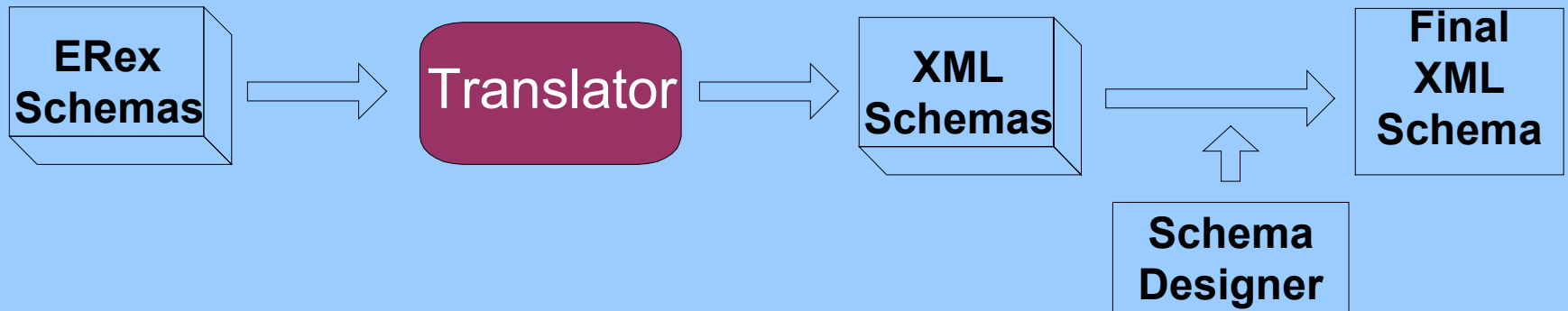
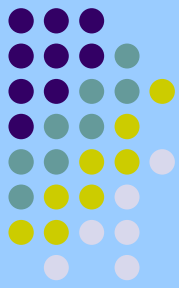
Translating ERex schemas to XML schemas

Oct 13, 2003

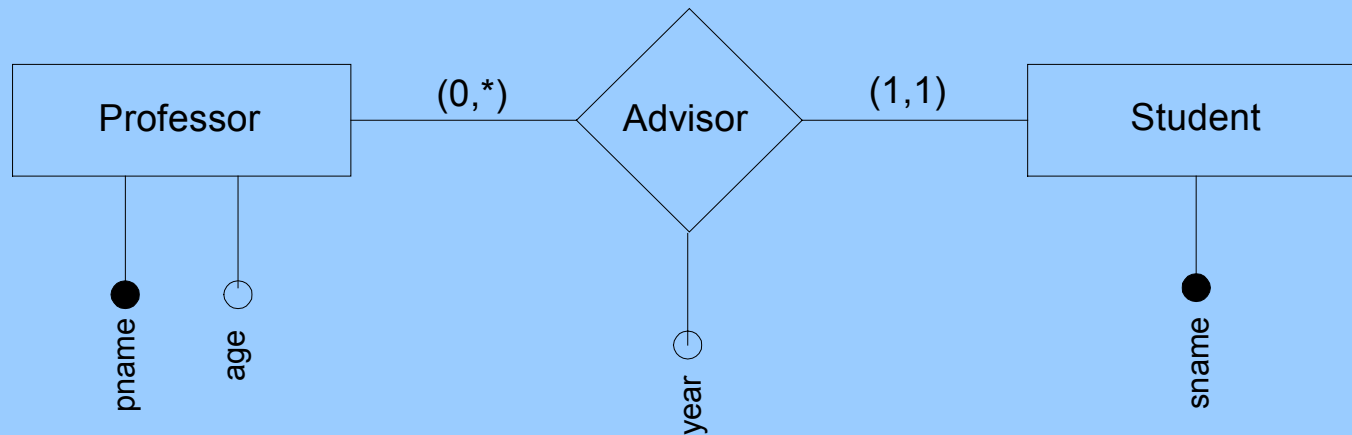
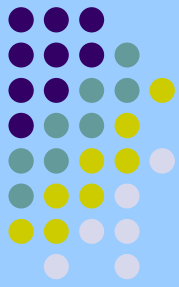
Murali Mani, Antonio Badia



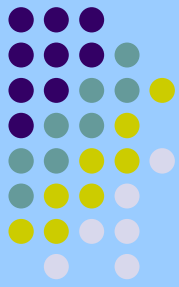
System Architecture



1:n relationships



Representing 1:n Relationships - subelement



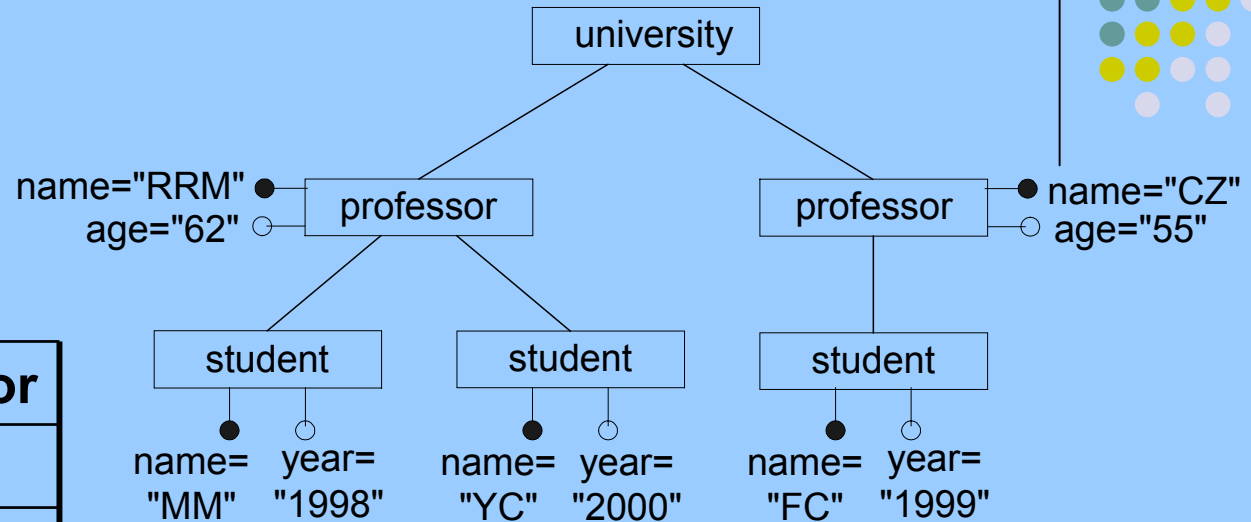
Professor

<u>name</u>	age
RRM	62
CZ	55

Student

<u>name</u>	year	advisor
MM	1998	RRM
YC	2000	RRM
FC	1999	CZ

Student (advisor) references
Professor (name)



University → university (Professor*)

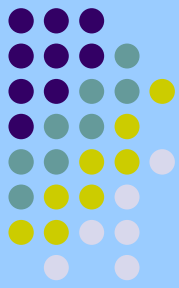
Professor → professor (@name, @age, Student*)

Student → student (@name, @year)

(University, Professor, <@name>)

(University, Student, <@name>)

Representing 1:n Relationships - IDREF



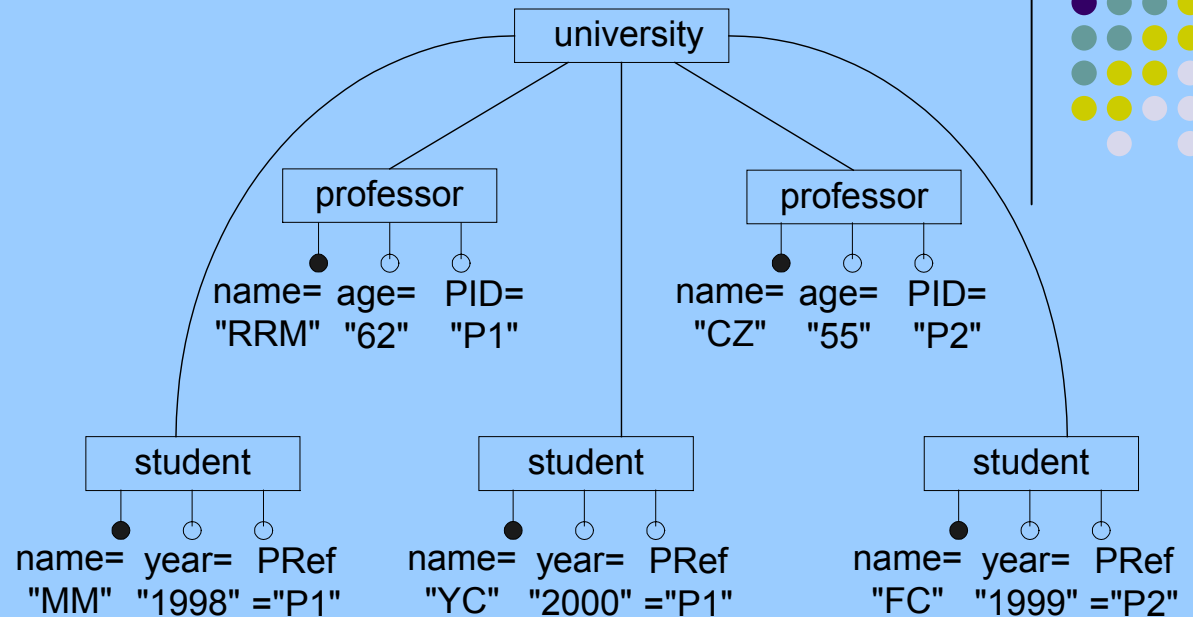
Professor

<u>name</u>	age
RRM	62
CZ	55

Student

<u>name</u>	year	advisor
MM	1998	RRM
YC	2000	RRM
FC	1999	CZ

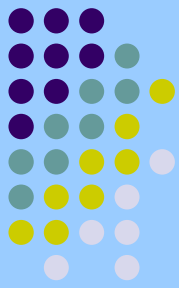
Student (advisor) references
Professor (name)



University → university (Professor*, Student*)
 Professor → professor (@name, @age, @PID)
 Student → student (@name, @year, @PRef)

(University, Professor, <@name>)
 (University, Student, <@name>)
@PRef::IDREF references (Professor)

Representing 1:n Relationships – foreign keys



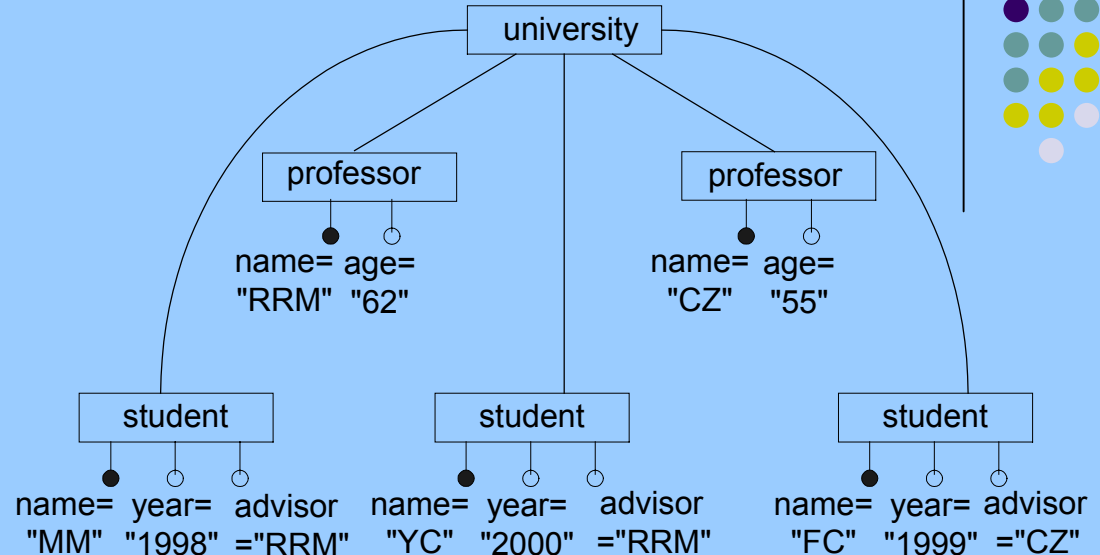
Professor

<u>name</u>	age
RRM	62
CZ	55

Student

<u>name</u>	year	advisor
MM	1998	RRM
YC	2000	RRM
FC	1999	CZ

Student (advisor) references
Professor (name)



University → university (Professor*, Student*)

Professor → professor (@name, @age)

Student → student (@name, @year, @advisor)

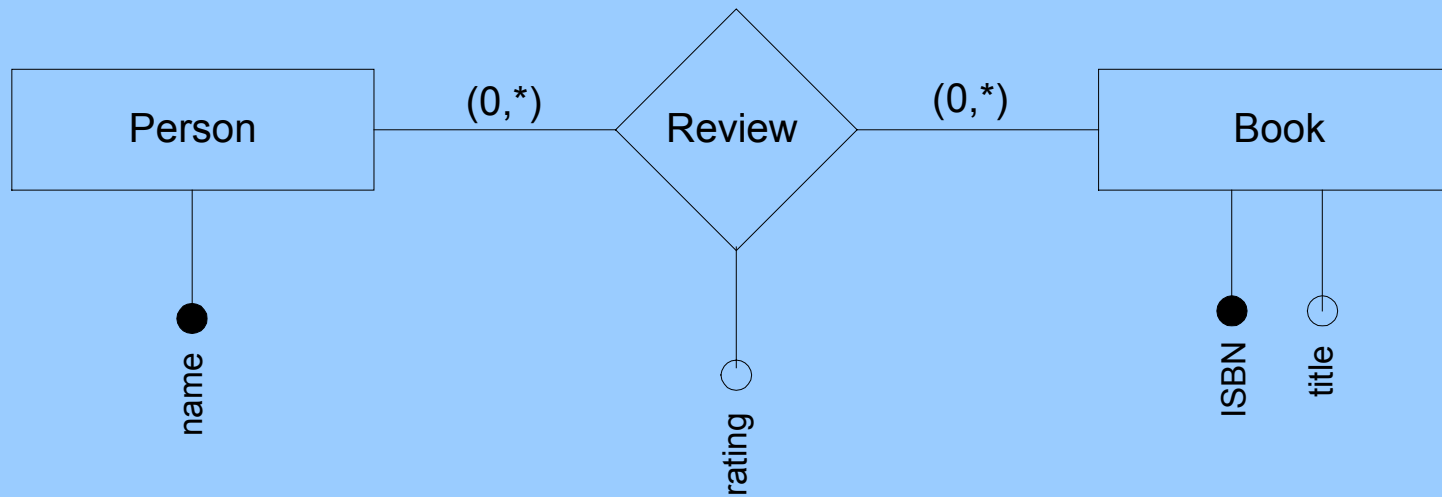
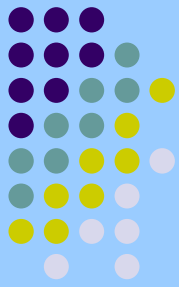
(University, Professor, <@name>)

(University, Student, <@name>)

(University, Student, <@advisor>) references

(University, Professor, <@name>)

m:n relationships



Representing m:n relationships

Person

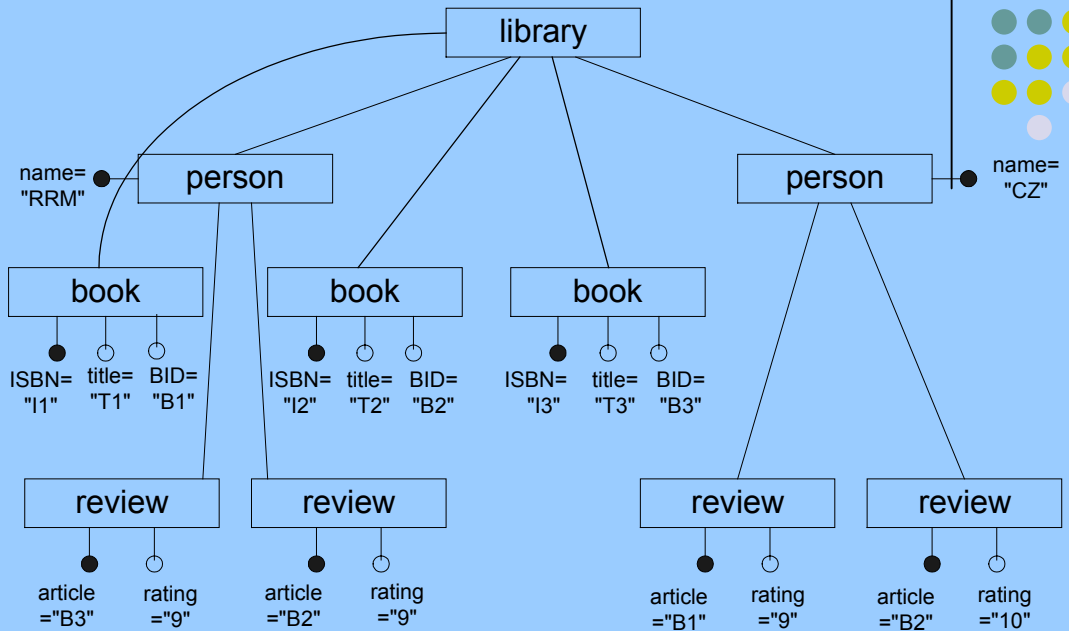
<u>name</u>
RRM
CZ

Book

<u>ISBN</u>	title
I1	T1
I2	T2
I3	T3

Review

<u>pname</u>	<u>ISBN</u>	rating
RRM	I3	9
RRM	I2	9
CZ	I1	9
CZ	I1	10



Library → library (Person*, Book*)
 Person → person (@name, Review*)
 Book → book (@ISBN, @title, @BID)
 Review → review (@article, @rating)

Review (pname) references
 Person (name)
 Review (ISBN) references
 Book (ISBN)

(Library, Person, <@name>) (Library, Book, <@ISBN>)
 (Person, Review, <@article>)
 @article::IDREF references (Book)

Representing m:n relationships

Person

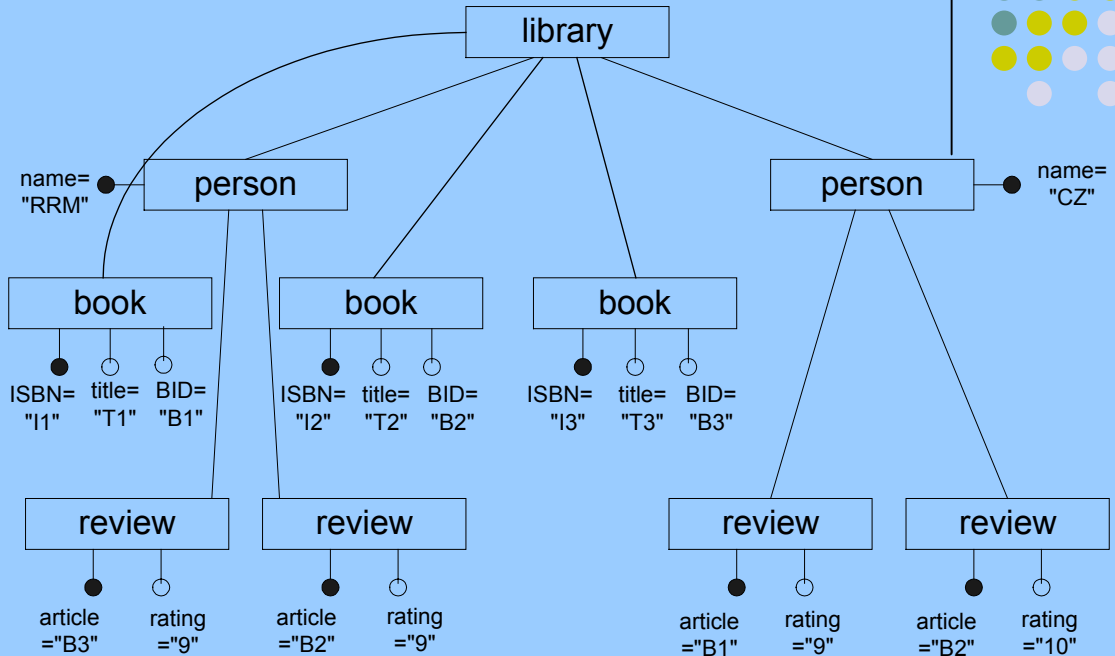
<u>name</u>
RRM
CZ

Book

<u>ISBN</u>	<u>title</u>
I1	T1
I2	T2
I3	T3

Review

<u>pname</u>	<u>ISBN</u>	rating
RRM	I3	9
RRM	I2	9
CZ	I1	9
CZ	I1	10

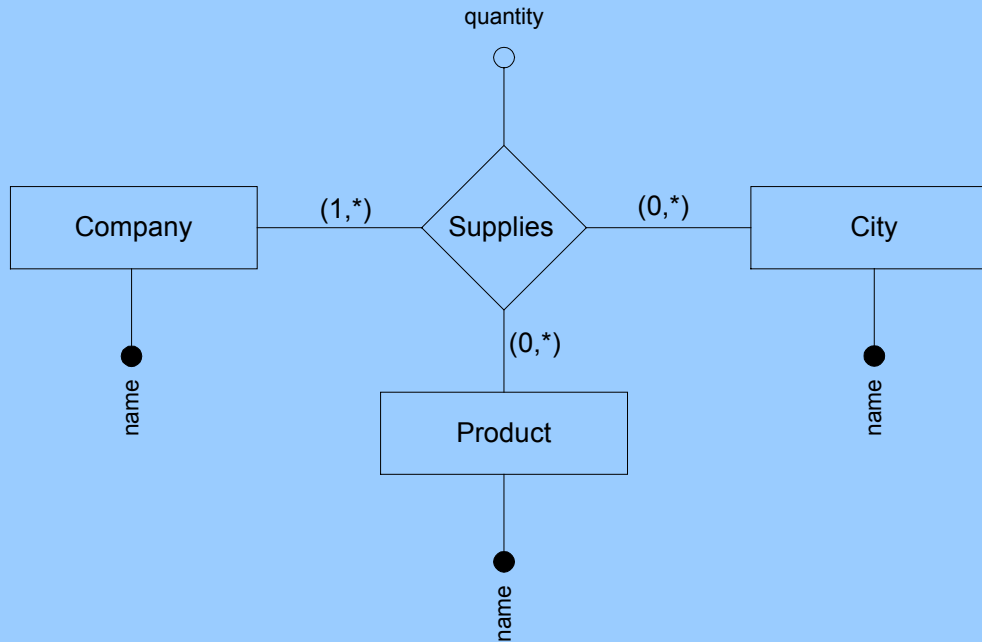
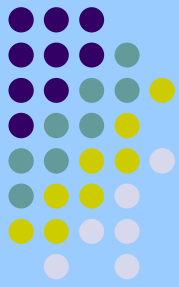


Query: What is the rating given by RRM for book T1

$$\pi_{\text{rating}} ((\sigma_{\text{title}=\text{T1}} (\text{Book})) \otimes (\sigma_{\text{pname}=\text{RRM}} (\text{Review})))$$

person[@name=RRM]/review[@article⇒Book/title=T1]/@rating

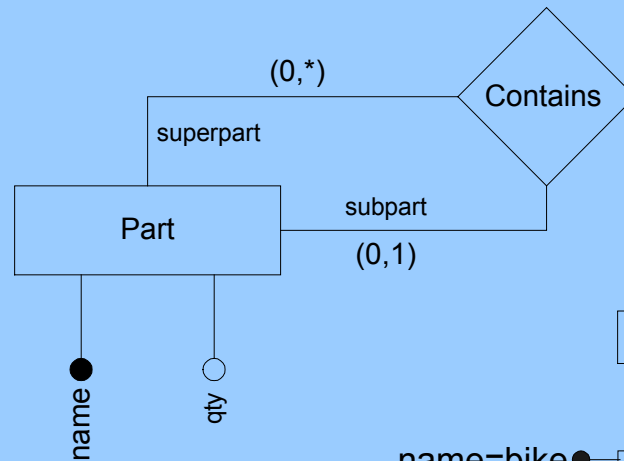
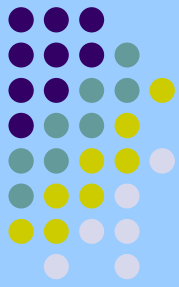
N-ary relationships



Root → root (Company*, Product*, City*)
Company → company (@name, Supply+)
Supply → supply (@ProdRef, @CityRef, @qty)
Product → product (@name, @ProdID)
City → city (@name, @CityID)

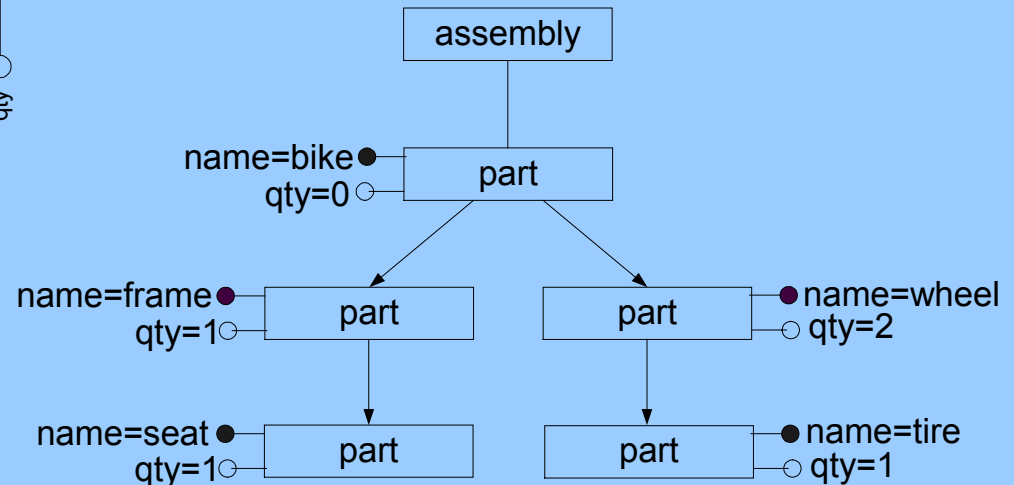
(Root, Company, <@name>)
(Root, Product, <@name>)
(Root, City, <@name>)
@ProdRef::IDREF references
(Product)
@CityRef::IDREF references
(City)

Recursive Relationships



Assembly

<u>name</u>	superPart	qty
seat	frame	1
tire	wheel	1
frame	bike	1
wheel	bike	2
bike	null	0

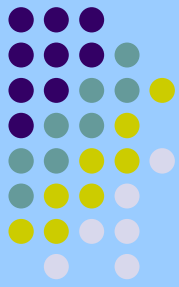


Assembly → assembly (Part*)

Part → part (@name, @qty, Part*)

(assembly, part, <@name>)

Ordered Relationships

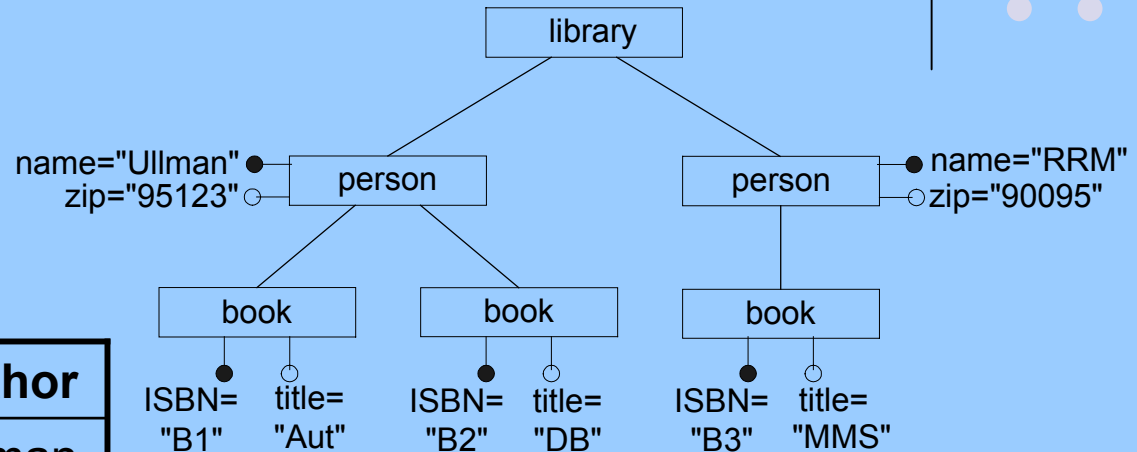


Person

<u>name</u>	zip
Ullman	95123
RRM	90095

Book

<u>ISBN</u>	title	order	author
B1	DB	2	Ullman
B2	Aut	1	Ullman
B3	MMS	1	RRM



Library → library (Person*)

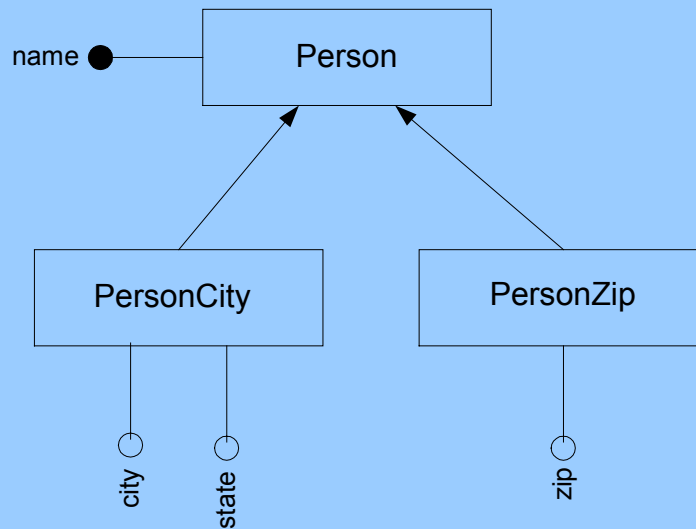
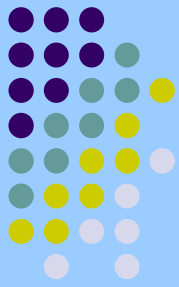
Person → person (@name, @zip, Book*)

Book → book (@ISBN, @title)

(Library, Person, <@name>)

(Library, Book, <@ISBN>)

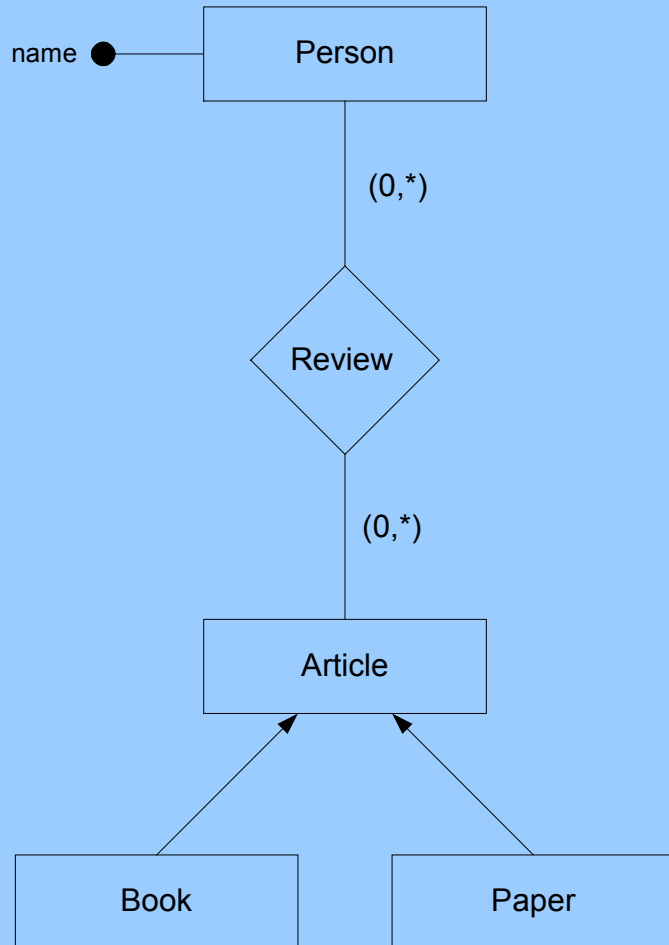
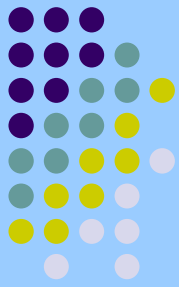
Categories and set constraints



$$\text{PersonCity} \cap \text{PersonZip} = \emptyset$$
$$\text{PersonCity} \cup \text{PersonZip} = \text{Person}$$

$\text{Person} \rightarrow \text{person} (@\text{name}, ((@\text{city}, @\text{state}) | @\text{zip}))$

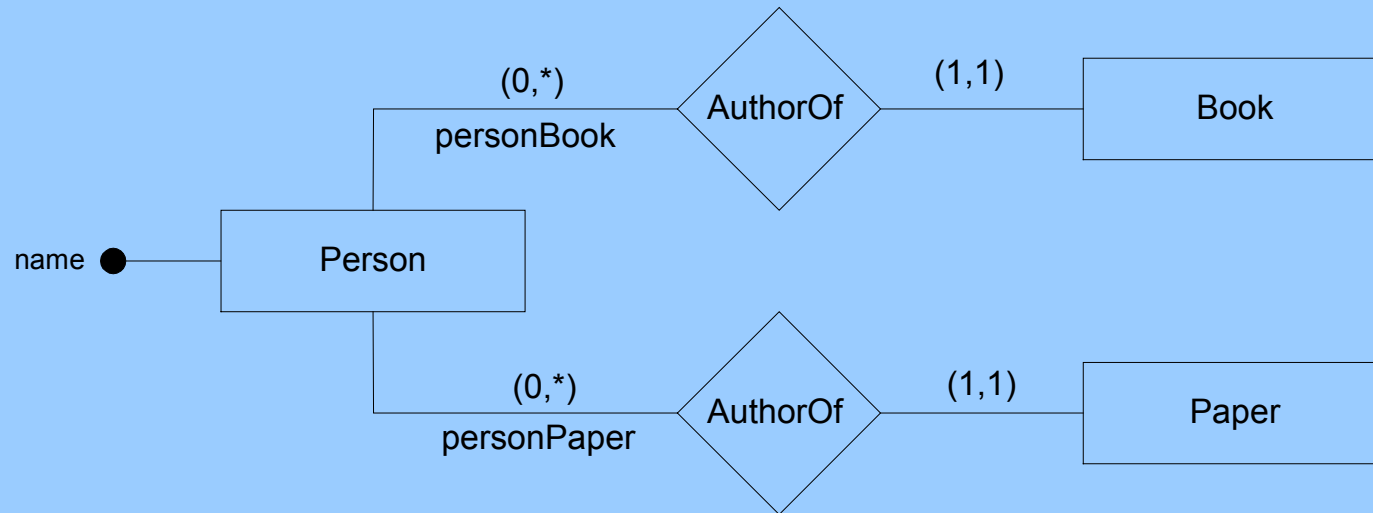
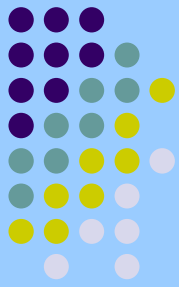
Categories and Set Constraints



Root \rightarrow root (Person*, Book*, Paper*)
Person \rightarrow person (@name, Review*)
Review \rightarrow review (@article, @rating)

(Root, Person, <@name>)
(Person, Review, <@article>)
@article::IDREF references (Book | Paper)

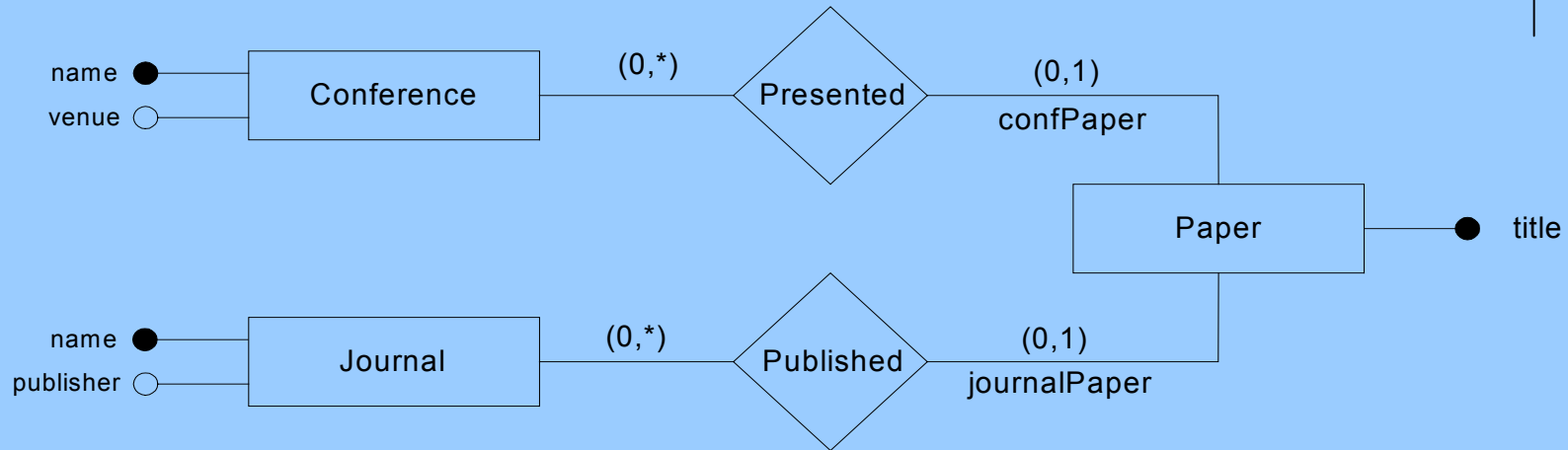
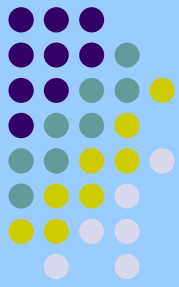
Set constraints on Roles



$$\text{personBook} \cap \text{personPaper} = \emptyset$$

$\text{Person} \rightarrow \text{person} (@\text{name}, (\text{Book}^* | \text{Paper}^*))$

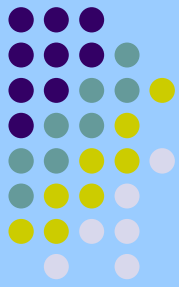
Set Constraints on Roles



$$\text{confPaper} \cap \text{journalPaper} = \emptyset$$
$$\text{confPaper} \cup \text{journalPaper} = \text{Paper}$$

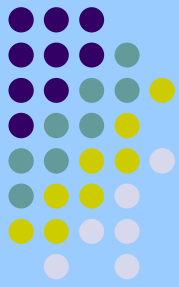
Conference \rightarrow conference (@name, @venue, Paper*)
Journal \rightarrow journal (@name, @publisher, Paper*)

Converting ERex → XML

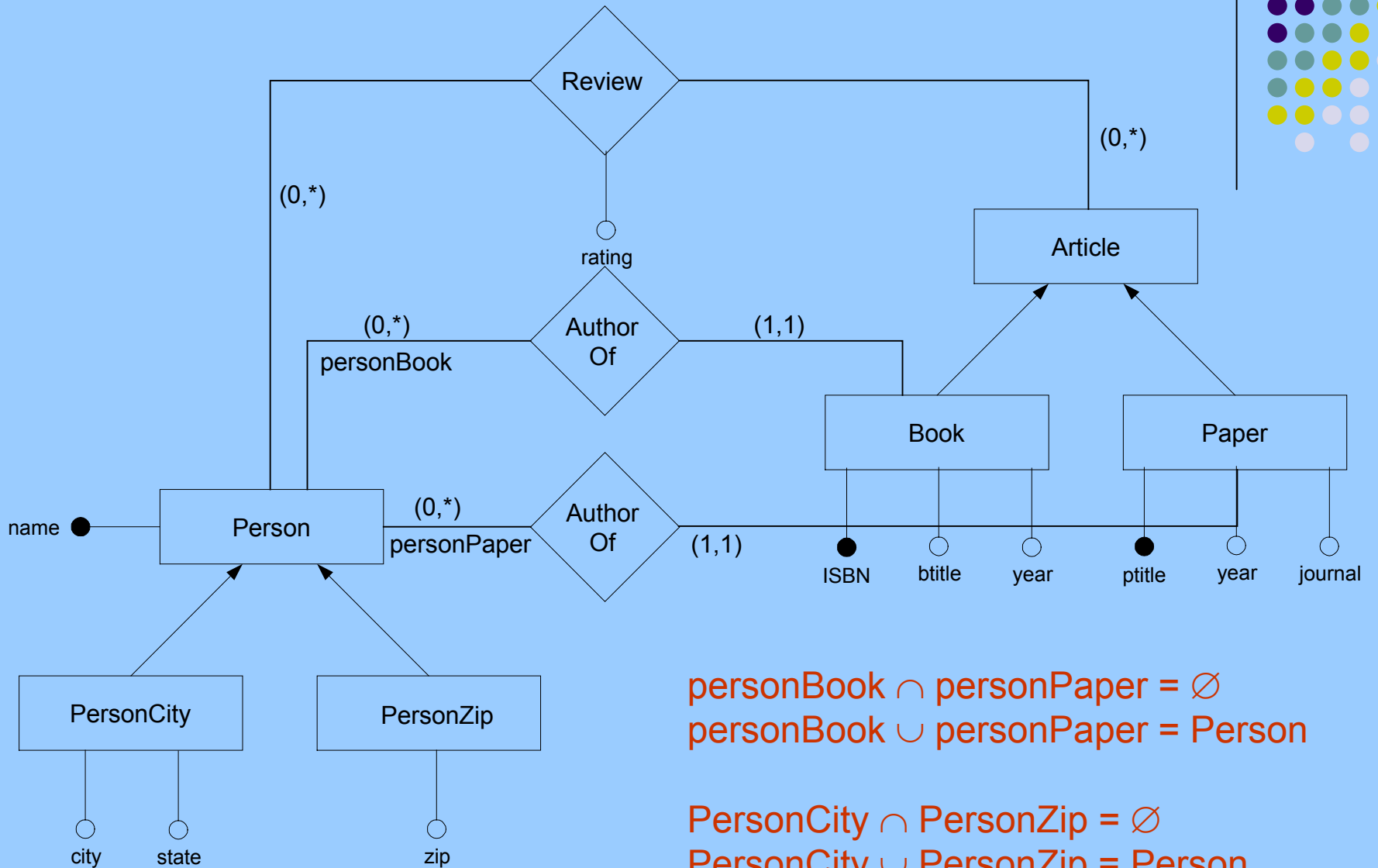
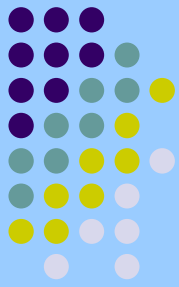


- Goals
 - Maximize relationships represented using subelement.
 - Others try to represent using IDREF

Algorithm: ERex \rightarrow XML



- A non-terminal symbol for each
 - entity type with key
 - m:n relationship
 - n-ary relationship
 - Root non-terminal symbol
- Represent attributes
- Represent relationships and identify top nodes
 - 1:1 and 1:n relationships
 - m:n relationships
 - n-ary relationships
- Identify key and IDREF constraints.

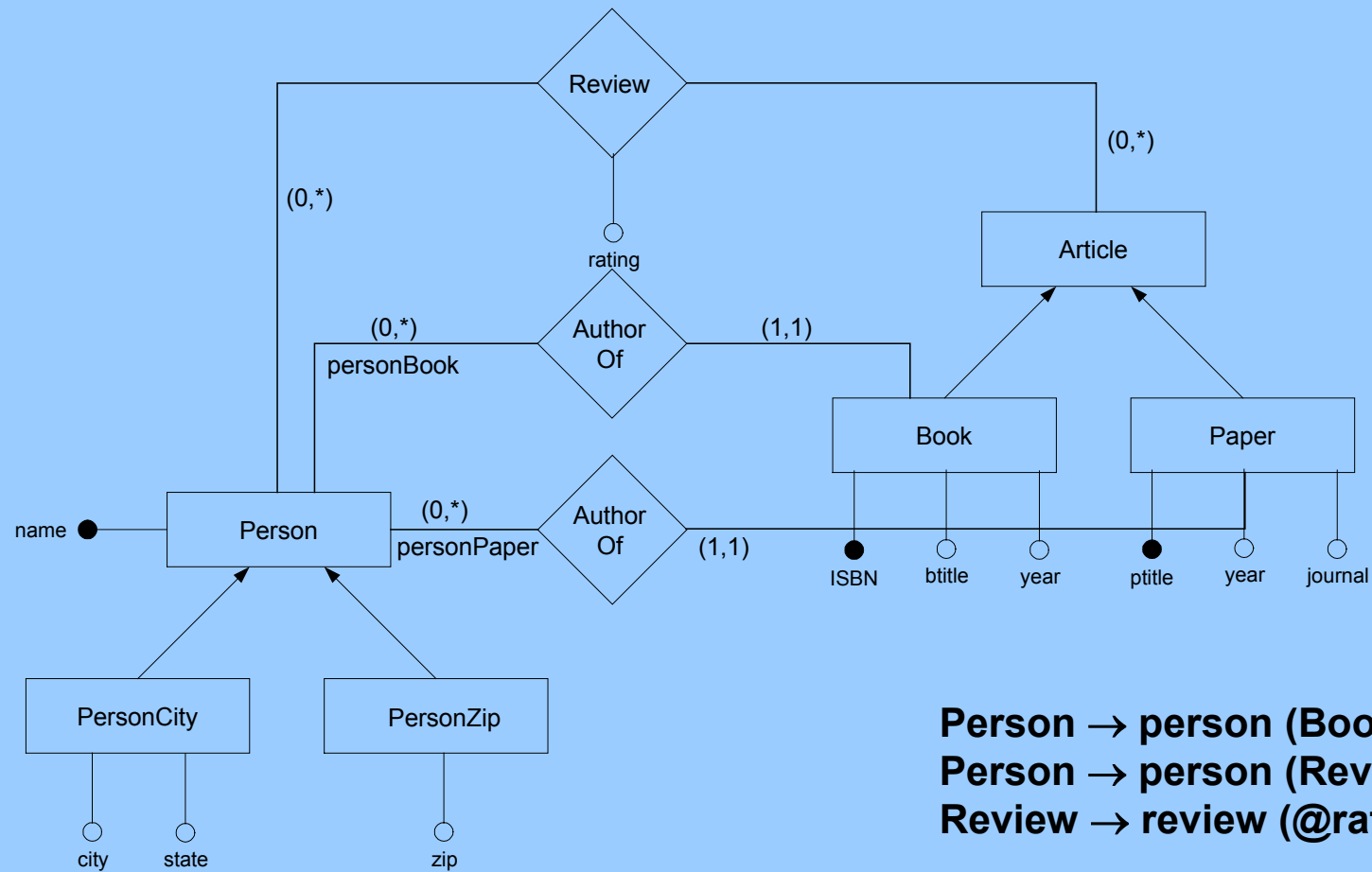
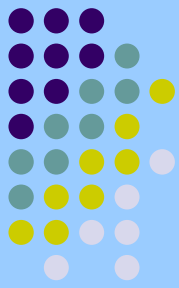


$$\text{personBook} \cap \text{personPaper} = \emptyset$$

$$\text{personBook} \cup \text{personPaper} = \text{Person}$$

$$\text{PersonCity} \cap \text{PersonZip} = \emptyset$$

$$\text{PersonCity} \cup \text{PersonZip} = \text{Person}$$



N = {Root, Person, Book, Paper, Review}

Book → book (@ISBN, @btitle, @year)

Paper → paper (@ptitle, @year, @journal)

Person → person (@name, ((@city, @state) | @zip))

Person → person (Book* | Paper*)

Person → person (Review*)

Review → review (@rating, @article)

Root → root (Person*)

(Root, Person, <@name>)

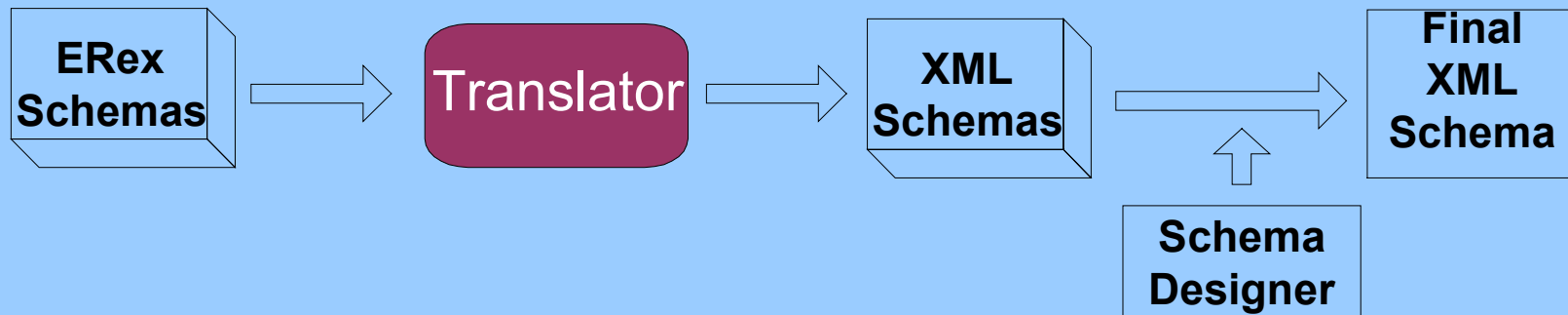
(Root, Book, <@ISBN>)

(Root, Paper, <@ptitle>)

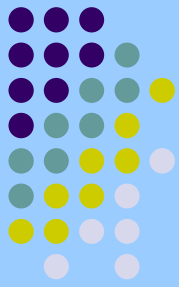
(Person, Review, <@article>)

@article::IDREF references (Book | Paper)

Conclusions



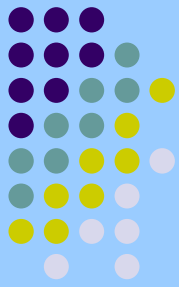
- Obtained good XML Schema from ERex schemas



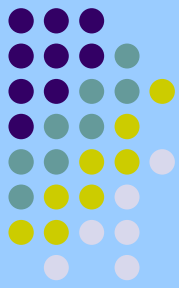
Part II:

Translation between Relational and XML models.

Why publish relational databases as XML?



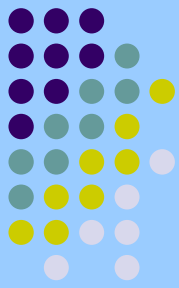
- Provide an XML view for our legacy data
 - Users/Applications can query our data over the web using standards.
 - Easier to query XML than legacy (relational) data?
- Convert our legacy data to XML
 - We can exchange data with applications.
 - Store data in XML databases?
 - Easier to query?



Application Scenarios

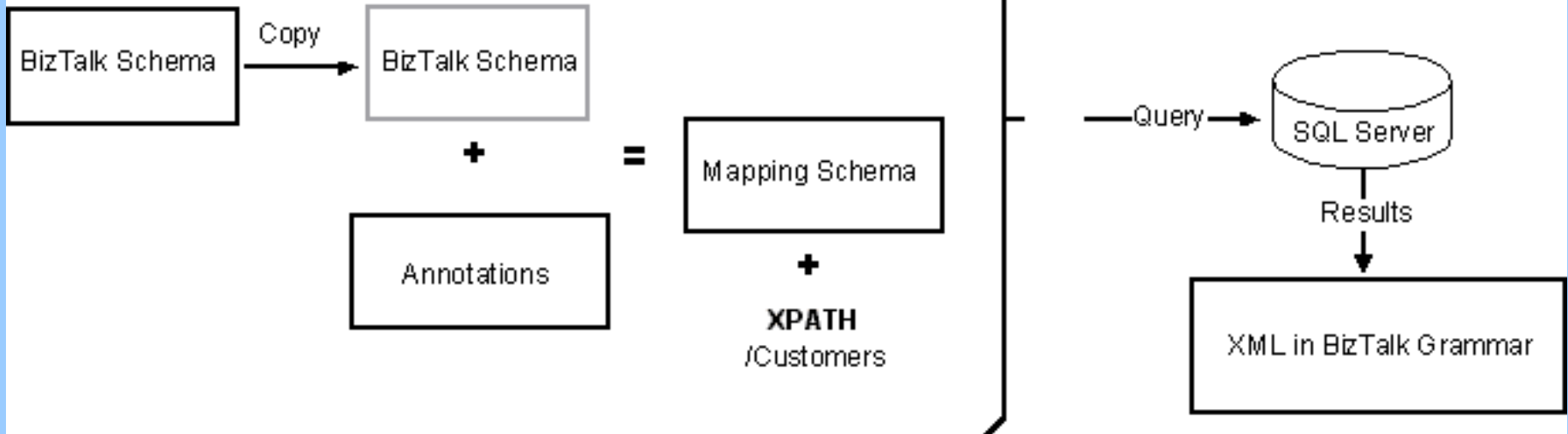
- Schema Matching Problem
 - Given a relational schema and an XML schema by a standards body, how do we map this relational schema to XML?
 - Tools such as XML Extender from IBM, Clio (University of Toronto and IBM), MS SQL Server
- Schema Mapping Problem
 - Given a relational schema, how do we come up with a **good** XML schema?

Schema Matching: MS SQL Server Architecture



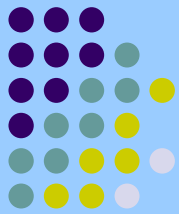
XML Views - BizTalk

Uses same schema format as BizTalk framework and server



COPYRIGHT NOTICE. Copyright © 2003 Microsoft Corporation, One Microsoft Way, Redmond, Washington 98052-6399 U.S.A. All rights reserved.

Schema Mapping Problem



Name	Age	Salary
A	25	20000
B	62	140000



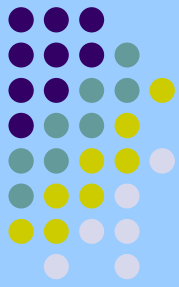
```
<person>
  <Name>A</name>
  <Age>25</Age>
  ...
</person>
```

- Users/Applications see a uniform XML view
- Exchange data with other applications
- Store XML in native XML databases
- Querying XML is easier?

Problem:

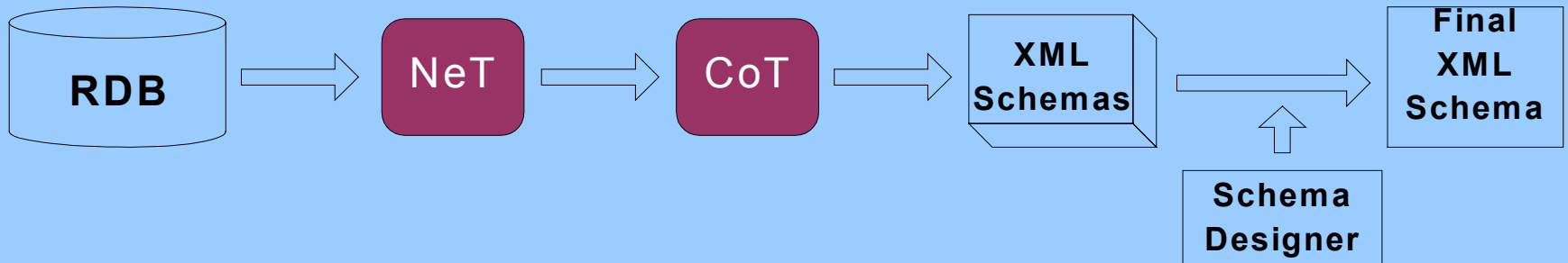
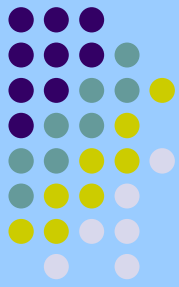
What is a good XML schema for a relational schema?

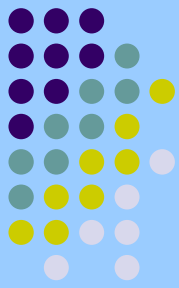
Goals



- XML Schema should maintain constraints.
- Resulting XML should not introduce redundancies.
- Most relationships can be navigated using path expressions, rather than joins.
- Minimal user interaction: Our translator should suggest good XML schemas to the database designer.

System Architecture

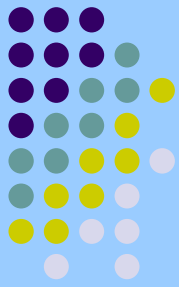




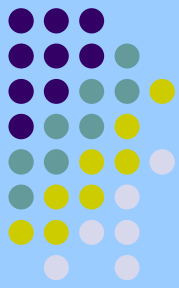
Related Work

- XML-DBMS
 - Template driven mapping language
- SilkRoute
 - Declarative Query Language (RXL) for viewing relational data as XML
- Xperanto
 - User specifies query in XML Query Language

Algorithms



- Naïve
 - FT (Flat Translation)
- Consider relational data
 - NeT (Nesting-based Translation)
- Consider relational schema
 - CoT (Constraint-based Translation)



FT: Flat Translation

- 1 : 1 mapping from relational to XML
- Idea
 - A type (non-terminal) corresponding to every relation
 - Attributes of a relation form attributes of the type
 - Keys and foreign keys are preserved

FT: Flat Translation - Example

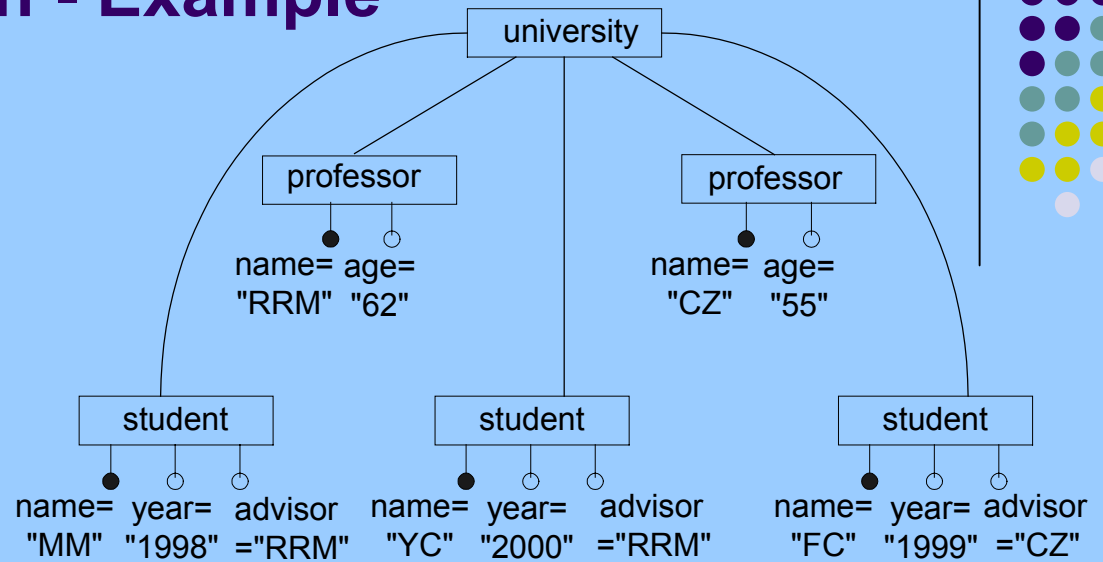
Professor

<u>name</u>	age
RRM	62
CZ	55

Student

<u>name</u>	year	advisor
MM	1998	RRM
YC	2000	RRM
FC	1999	CZ

Student (advisor) references
Professor (name)



University → university (Professor*, Student*)

Professor → professor (@name, @age)

Student → student (@name, @year, @advisor)

(University, Professor, <@name>)

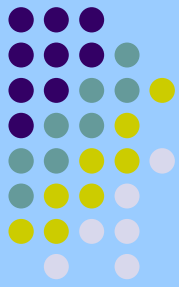
(University, Student, <@name>)

(University, Student, <@advisor>) references

(University, Professor, <@name>)



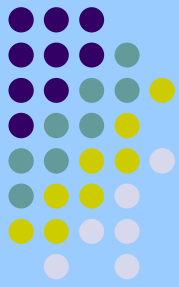
NeT: Nesting-based Translation



- Idea:

Make use of non-flat features provided by XML: represent repeating groups using *, +

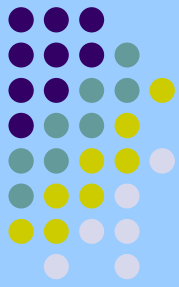
NeT: Example



Course (cname, prof, text)

cname	prof	text
Algorithms	Gafni	Udi Manber
Algorithms	Gafni	CLR
Algorithms	Majid	Udi Manber
Algorithms	Majid	CLR

NeT: Example



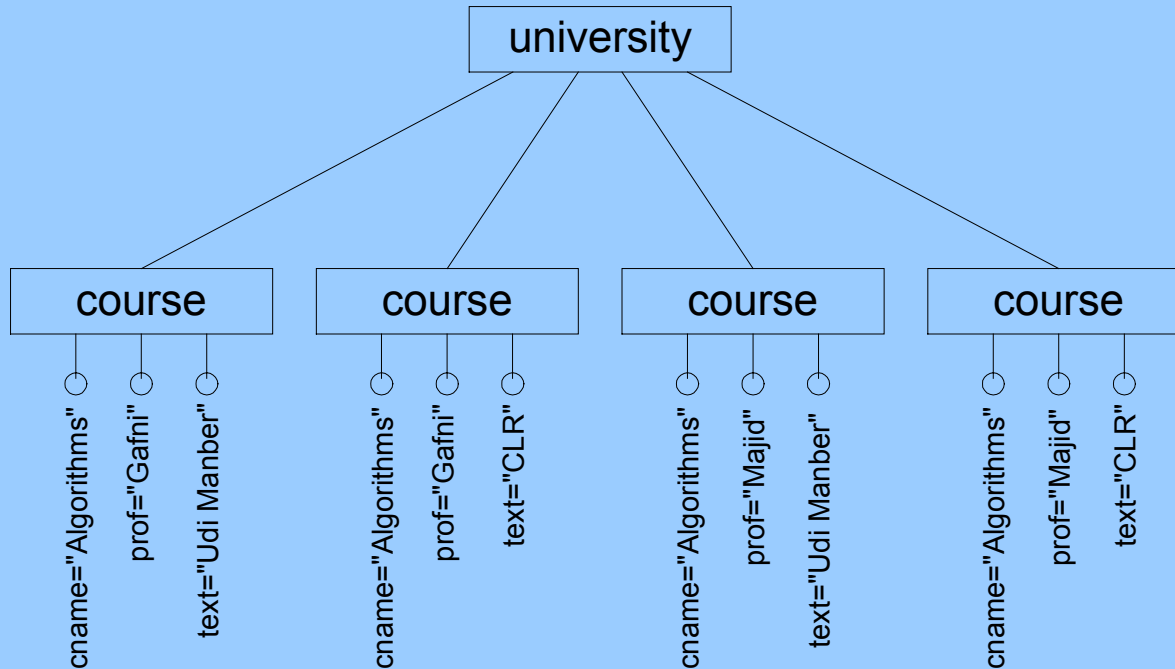
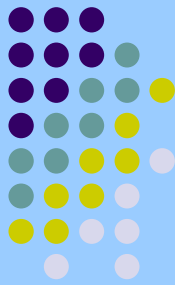
Course (cname, prof[±], text)

cname	prof	text
Algorithms	{Gafni, Majid}	Udi Manber
Algorithms	{Gafni, Majid}	CLR

Course (cname, prof[±], text[±])

Cname	prof	text
Algorithms	{Gafni, Majid}	{Udi Manber, CLR}

NeT: Example

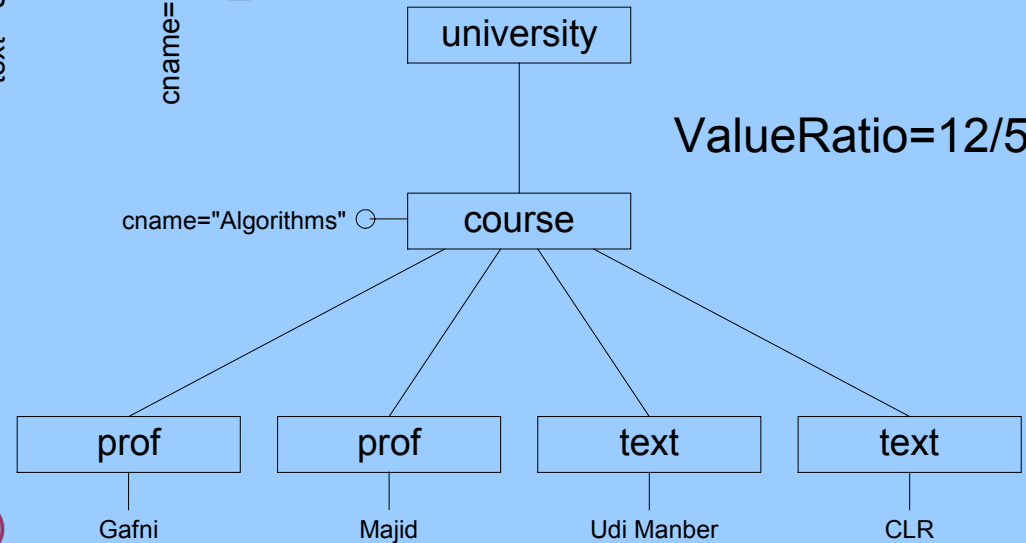


University \rightarrow university (Course*)
 Course \rightarrow course (@cname, @prof, @text)

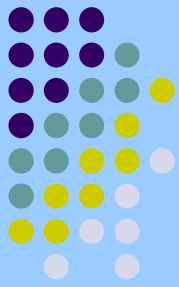
(University, Course, <@cname, @prof, @text>)

University \rightarrow university (Course*)
 Course \rightarrow course (@cname, Prof+, Text+)
 Prof \rightarrow prof (#PCDATA)
 Text \rightarrow text (#PCDATA)
 #PCDATA \rightarrow pcdta (ϵ)

(University, Course, <@cname, prof, text>)



ValueRatio=12/5



NeT: Example

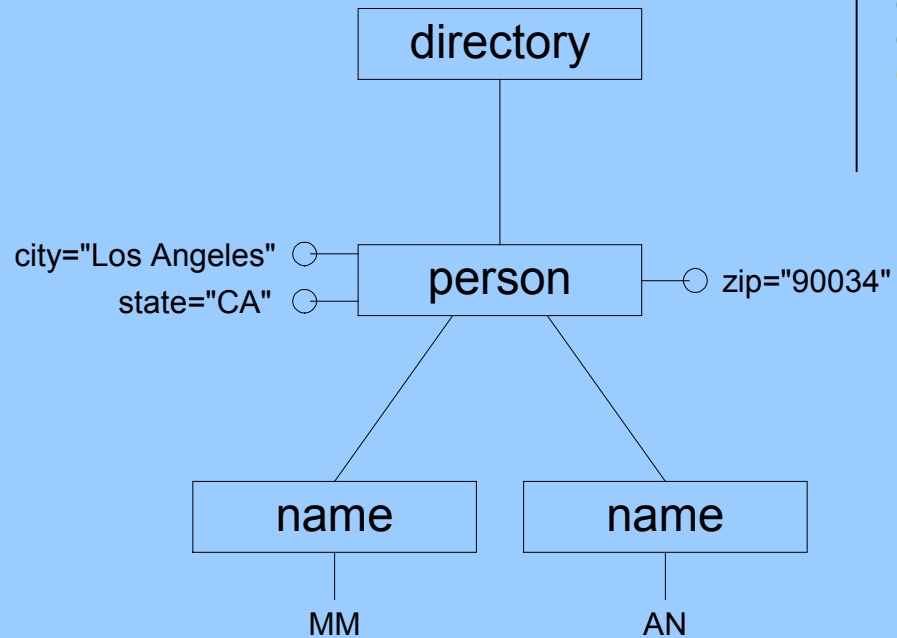
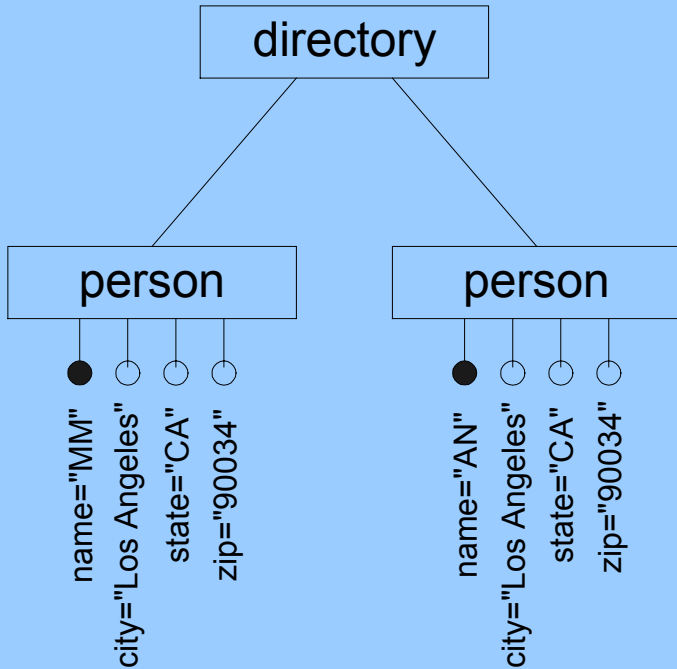
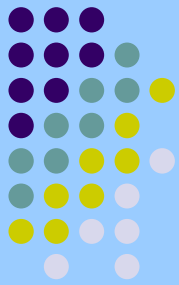
Person (name, city, state, zip)

name	city	state	zip
MM	Los Angeles	CA	90034
AN	Los Angeles	CA	90034

Person (name⁺, city, state, zip)

{MM, AN}	Los Angeles	CA	90034
----------	-------------	----	-------

NeT: Example



Directory → directory (Person*)
 Person → person (@name, @city,
 @state, @zip)

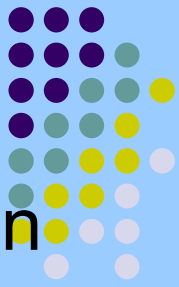
(directory, person, <@name>)

Directory → directory (Person*)
 Person → person (@city, @state,
 @zip, Name*)
 Name → name (#PCDATA)

(directory, person, <@name>)

ValueRatio=8/5

NeT: Summary

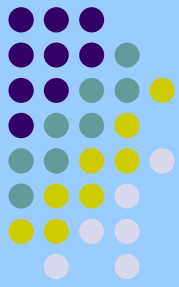


- Consider Table t with column set C . Nesting on column X is defined as:

Any two tuples with the same values for $(C - X)$ will be combined to one tuple

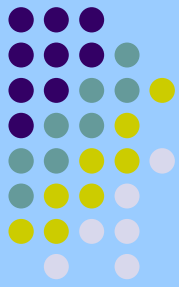
- **Observation: We need to nest only on key columns**
- Advantages of NeT
 - NeT removes redundancy if relation is not in 4NF
 - NeT provides more intuitive XML schemas with less redundancy

NeT: Experimentation



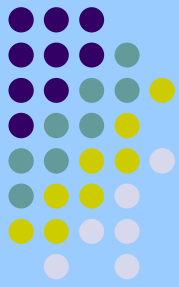
Test Set	#attr	#tuples	ValueRatio	#nested attributes	Time (sec)
Balloons1	5	16	3.64	3	1.08
Hayes	6	132	1.52	1	1.01
Bupa	7	345	1	0	4.40
Balance	5	625	2.79	4	21.48
TA_Eval	6	110	1.24	5	24.83
Car	7	1728	15.53	6	469.47
Flare	13	365	1.67	4	6693.41

CoT: Constraint-based Translation



- Translating relational schema
- Idea:
Use foreign key constraints and our knowledge of how to represent relationships to come up with “better” XML models.

CoT: Step 1



Professor

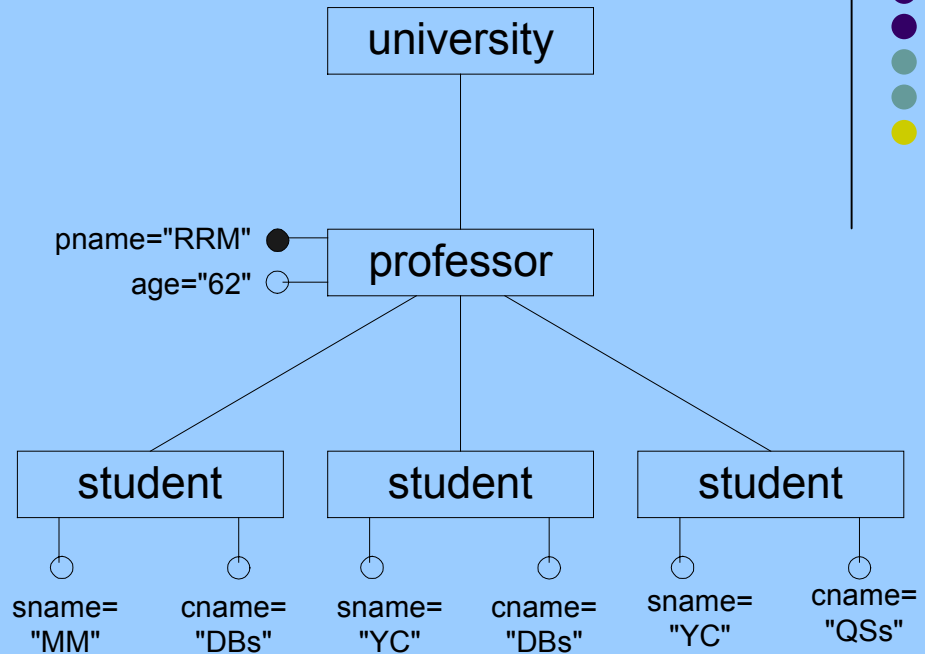
<u>pname</u>	age
RRM	62

Student

<u>sname</u>	<u>advisor</u>	<u>cname</u>
MM	RRM	DBs
YC	RRM	DBs
YC	RRM	QSs

Student (advisor) references
Professor (pname)

ValueRatio=11/8



University → university (Professor*)

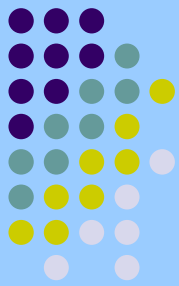
Professor → **professor(@pname, @age, Student*)**

Student → student (@sname, @cname)

(University, Professor, <@pname>)

(Professor, Student, <@sname, @cname>)

CoT: Step 2



Professor

<u>pname</u>	age
RRM	62

Course

<u>cname</u>	since
DBs	1979
QSS	1962

Student

<u>sname</u>	<u>advisor</u>	<u>cname</u>
MM	RRM	DBs
YC	RRM	DBs
YC	RRM	QSS

Student (advisor) references
 Professor (pname)
 Student (cname) references
 Course (cname)

University → university (Professor*, Course*)

Professor → professor (@pname, @age, Student*)

Student → student (@sname, @CRef)

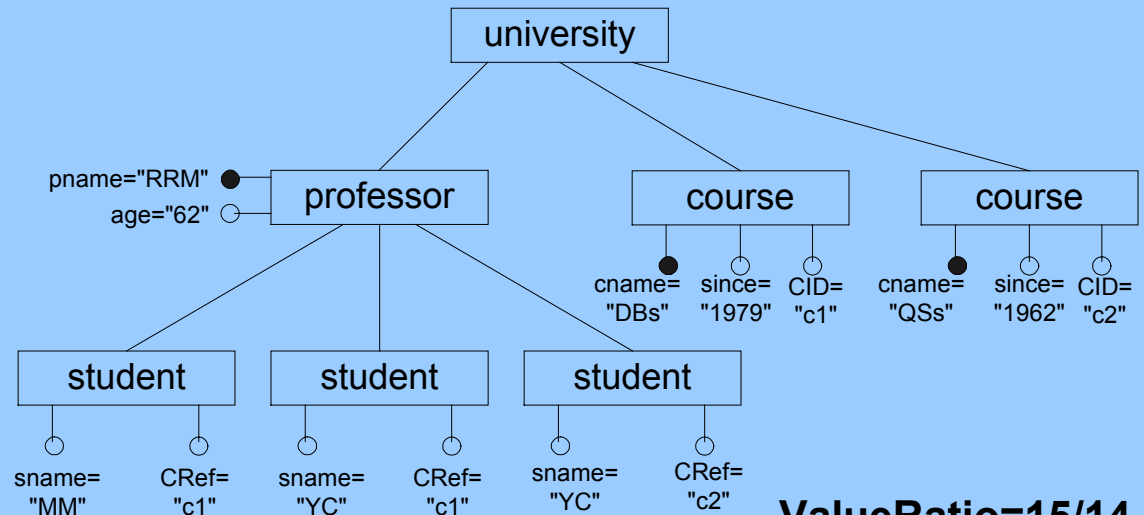
Course → course (@cname, @since, @CID)

(University, Professor, <@pname>)

(University, Course, <@cname>)

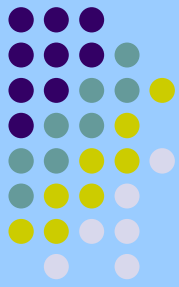
(Professor, Student, <@sname, @Cref>)

CRef::IDREF references (Course)



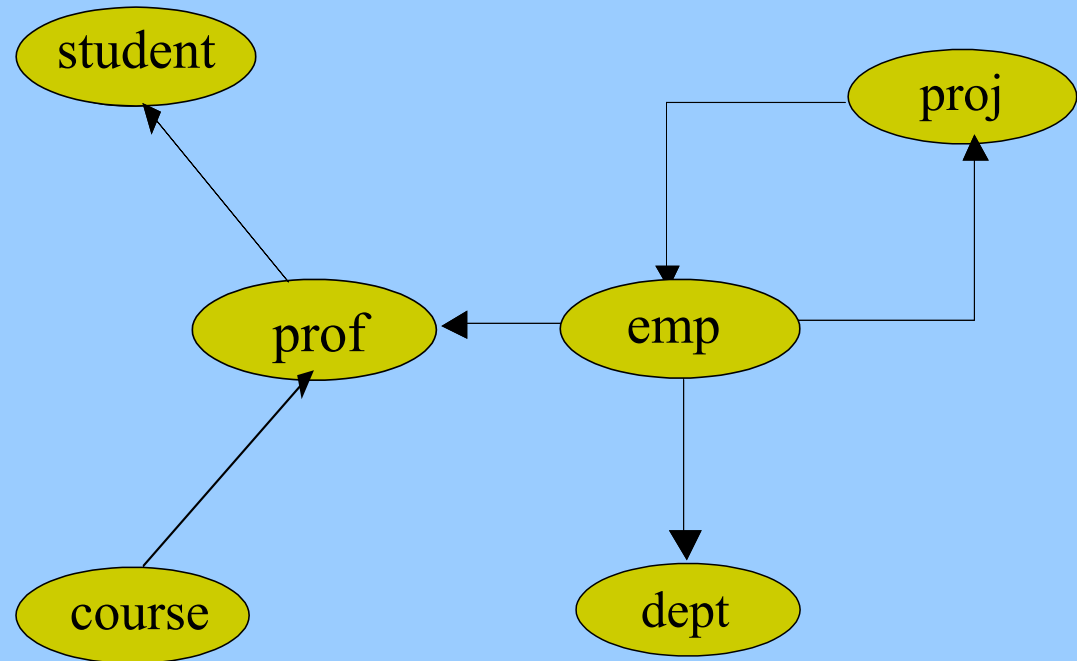
ValueRatio=15/14

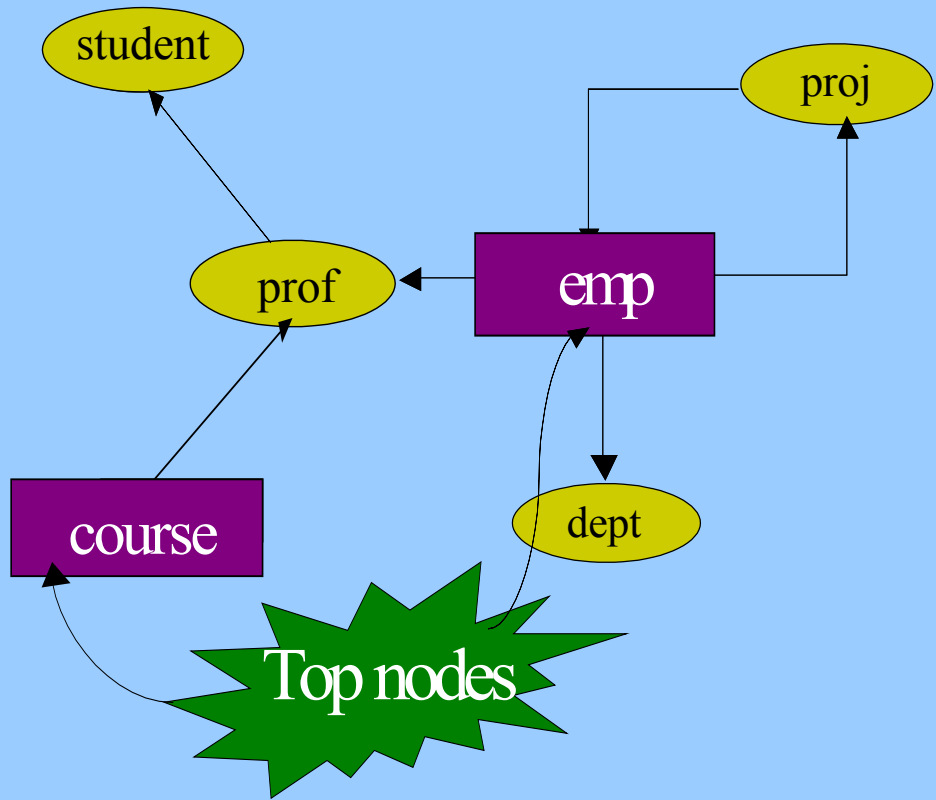
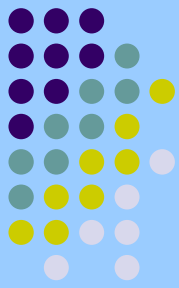
CoT: Example



Student (SID, name, advisor)
Emp (EID, name, projName)
Prof (EID, name, teach)
Course (CID, title, room)
Dept (dno, mgr)
Proj (pname, pmgr)

Student (advisor) references
Prof (EID)
Emp (projName) references
Proj (pname)
Prof (teach) references
Course (CID)
Prof (EID, name) references
Emp (EID, name)
Dept (mgr) references
Emp (EID)
Proj (pmgr) references
Emp (EID)





(University, Course, <@CID>)
(University, Emp, <@EID>)
(University, Prof, <@EmpRef>)
(University, Student, <@SID>)
(University, Proj, @pname)
(University, Dept, @dno)

**@EmpRef::IDREF references
(Emp)**

**@ProjRef::IDREF references
(Proj)**

University → university (Course*, Emp*)

Course → course (@CID, @title, @room, Prof*)

Prof → prof (@EmpRef, Student*)

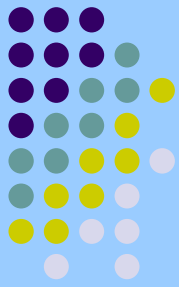
Student → student (@SID, @name)

Emp → emp (@EID, @name, @ProjRef, @EmpID, Dept*, Proj*)

Proj → proj (@pname, @ProjID)

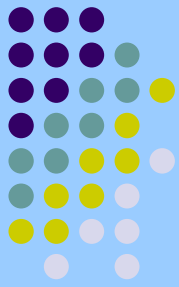
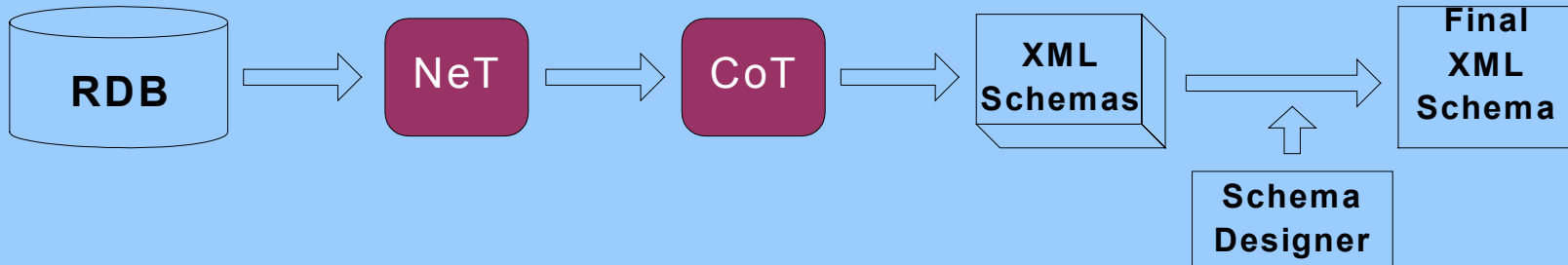
Dept → dept (@dno)

CoT Experimentation

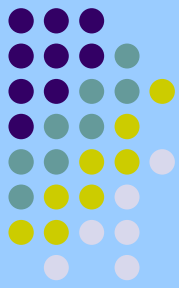


- Ran on TPC-H data
- Value ratio $> 100/88$ (size decreased by more than 12%)

Conclusions

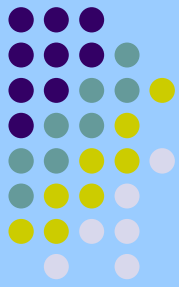


- We obtained “good” XML schemas from relational schemas
- Constraints are maintained
- Redundancies are decreased
- Most relationships can be navigated using path expressions.
- Minimum user interaction



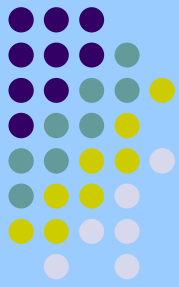
Storing XML data in relational databases

Options for storing XML data



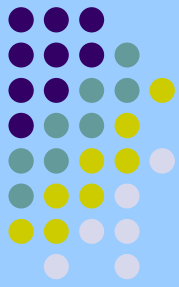
- Store in relational databases
 - Relational databases are robust and efficient (IBM XML Extender, Oracle, MS SQL Server)
- Store in native XML databases
 - More efficient for XML than relational databases (Natix, eXist, Tamino)
- Store in a combination of both
 - Structured portion of XML data in relational database, and unstructured portion in native XML store.

Related Work

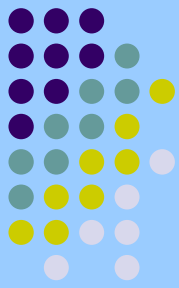


- STORED
 - No Schema
 - Use data mining techniques to find structured and frequent patterns
 - These are stored in relational DB, others in semi-structured overflow store
 - Drawback: Requires integration of relational DB and semi-structured store
- Storing paths
 - One relation for storing nodes, one for storing edges.
 - Drawback: Type information is lost.

Type-based relational storage



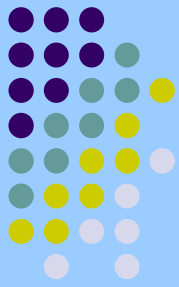
- Jayavel Shanmughasundaram
 - Several key ideas such as schema simplification, inlining, handling recursion
- LegoDB
 - Use the query workload to come up with an “efficient” relational schema.



Main features

- The entire XML document is shredded and stored in a relational database.
- All semantic constraints in XML schema are not captured in relational schema.
- We do not discuss how operations on XML are translated to SQL.

Why not capture all constraints?



$A \rightarrow a (@b, ((@c, D)^* | (@e, F)^*))$

(Root, A, <@b>)

A (@b, @c, @e)

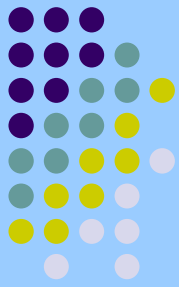
D (@aRef)

F (@aRef)

Semantic constraints lost

- if D refers to an A, then the corresponding @c should be non-null
- if F refers to an A, then the corresponding @e should be non-null

NF2 representation of regular tree grammars



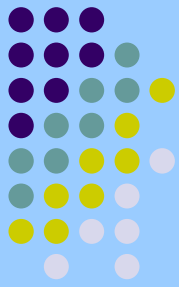
Eliminate union, “|”

paper \rightarrow ($@$ ptitle, $@$ journal | $@$ conference)

paper \rightarrow ($@$ ptitle, $@$ journal)

paper \rightarrow ($@$ ptitle, $@$ conference) **(or)**

paper \rightarrow ($@$ ptitle, $@$ journal, $@$ conference)

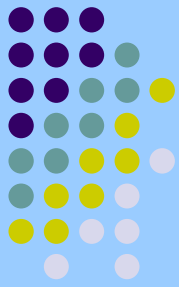


Schema Simplification

$$A \rightarrow a (@d, (B, C, B)^*)$$
$$A \rightarrow a (@d [1, 1], (B [2, 2], C [1, 1]) [0, *])$$
$$A \rightarrow a (@d [1, 1], B [2, 2] [0, *], C [1, 1] [0, *])$$
$$A \rightarrow a (@d [1, 1], B [0, *], C [0, *])$$
$$A \rightarrow a (@d, B^*, C^*)$$

Semantic information lost

- The number of B's is two times the number of C's
- for every C, there is a B that occurs before it, and one that occurs after it.



Inlining

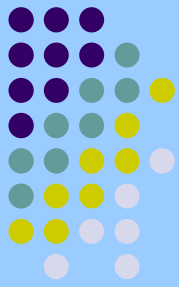
Conf → conf (@ctitle, @date, Venue)

Venue → venue (@city, @country)

Conf → conf (@ctitle, @date, @city, @country)

Why inlining?

- Lesser joins, hence more efficient



Mapping Collection Types

Conf → conf (@ctitle, @date, Paper*)

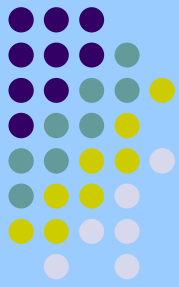
Paper → paper (@ptitle, @author)

Conf (@ctitle, @date)

Paper (@ptitle, @author, @confRef)

- Separate relation with foreign key for every collection type

IDREF attribute



Person → person (@name, @zip, Review*)

Book → book (@ISBN, @btitle, @BID)

Paper → paper (@ptitle, @journal, @PID)

Review → review (@article, @rating)

@article::IDREF references (Book | Paper)

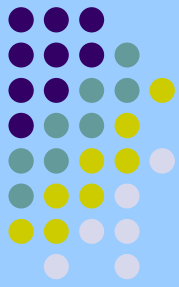
Person (@name, @zip)

Book (@ISBN, @btitle, @BID)

Paper (@ptitle, @journal, @PID)

Review (@personRef, @bookRef, @paperRef, @rating)

Recursion using ?



$A \rightarrow a (@d, A?)$

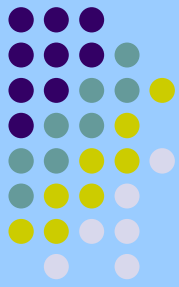
$A (@d, ARef)$

ARef refers to the child of A, and can be null

$A \rightarrow a (@d, ARef)$

ARef refers to the parent of A and can be null.

Recursion using *

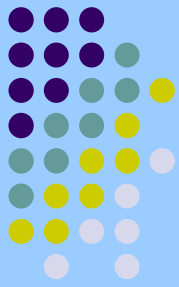


$A \rightarrow a (@d, A^*)$

$A (@d, ARef)$

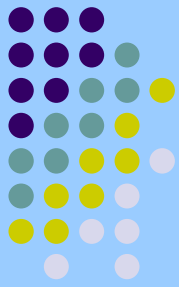
ARef refers to the parent of A, and can be null

Recursion – General technique



- For every cycle, we must have a separate relation
- Algorithm
 - For every strongly connected component, define a separate relation for one of the types.
 - In a strongly connected component, if there is a type which can be children of multiple types, then define a separate relation for that type.

Capturing Order in the Document



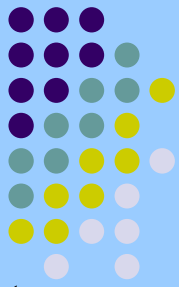
- Through order attributes
- Corresponding to each type in XML schema, say A, we have an associated order attribute, say aOrder

Conf → conf (@ctitle, @date, Venue)

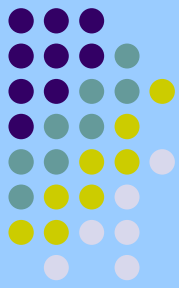
Venue → venue (@city, @country)

Conf → conf (@ctitle, @date, @confOrder, @city, @country, @venueOrder)

Conclusions

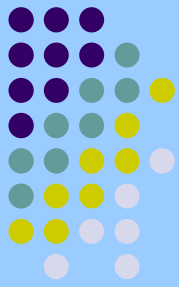


- XML schema with no recursion can be translated to relational schema with no nulls.
- XML schema with recursion cannot be translated to relational schema with no nulls.
- If recursion, separate relation needed for every cycle.
- All semantic constraints in XML cannot be captured in relational schema.
- XML resulting from CoT can be translated to the original relational schema; all semantic constraints are maintained.



Open Problems

- Standards specification: What structural and constraint specification schemes for XML are needed for database applications?
- XML used for text/document publishing: Keyword Search in XML documents
- Storing data consisting of structured and unstructured portions: integrating relational and XML stores.



Open Problem (contd...)

- Translating operations in XML model to underlying sources (relational)
 - Use annotated schema (MS SQL Server)
 - Use implicit annotations (LegoDB)
 - Query minimization: When we do automatic translation, we might perform unnecessary joins?