

# A Normal Form for XML Documents

---

- Overview of Relational Database Design Process
  - Functional Dependencies and Normalization
    - functional dependencies (FDs)
    - redundancy and update anomalies
    - third normal form (3NF) and Boyce-Codd normal form (BCNF)
    - design algorithms for 3NF and BCNF
  - Nested Normal Form for nested relations
- Normal Form for XML documents
  - redundancy and update anomalies for XML documents
  - functional dependencies
  - XNF: a normal form for XML documents
  - a design algorithm for XNF

This section is based on the paper A Normal Form for XML Documents by M. Arenas L. Libkin in Proceedings of ACM PODS02.

# A motivation Example for Normal Form Relations

<b>course</b>	<b>title</b>	<b>student_id</b>	<b>Name</b>	<b>Major</b>	<b>Grade</b>
391	Database	1234	Sarah	CS	9
391	Database	4321	Tom	CS	8
391	Database	2345	Bill	CS	7
201	Program	1234	Sarah	CS	6
201	Program	2345	Bill	CS	5

## Motivation Example

StudentCourse = ( course, title, student\_id, name, major, grade)

Student = ( student\_id, name, major)

Course = ( Course, title )

Registration = ( course, student\_id, grade)

# Functional Dependencies

## ■ Functional dependencies (FDs)

Let  $R$  be a relation scheme, and  $X \subseteq R$  and  $Y \subseteq R$  be sets of attributes.

Then the functional dependency

$$X \rightarrow Y$$

holds on  $R$  if in any legal relation  $r$ , for all pairs of tuples  $t_1$  and  $t_2$  in  $r$

$$t_1[X] = t_2[X] \quad \Rightarrow \quad t_1[Y] = t_2[Y].$$

Example:  $\text{student\_id} \rightarrow \text{name}$   
 $\text{course, student\_id} \twoheadrightarrow \text{grade}$

## Desirable Properties of Decomposition

- Minimizing redundancy
  - Boyce-Codd normal form
  - third normal form

## Boyce-Codd Normal Form

- A relation scheme  $R$  is said to be in Boyce-Codd normal form (BCNF) if for any non-trivial FD  $X \rightarrow A$  which holds in  $R$ ,  $X$  is a key of  $R$ , that is,  $X \rightarrow A$  holds in  $R$ .
  - no partial redundancy
  - no transitive redundancy
- Let  $U$  be a set of attributes,  $F$  be a set of FDs, and  $D = \{R_1, \dots, R_n\}$  be a decomposition of  $U$ . Then  $D$  is said to be a BCNF decomposition of  $U$  with respect to  $F$  if
  - $R$  is a join loss-less decomposition of  $U$  wrt  $F$ , and
  - every relation scheme  $R_i$  in  $D$  is in BCNF wrt  $F$ .

■ Example

Regist (Course#, Student#, Grade, Address, Phone)

is not in BCNF since

Student#  $\rightarrow$  Address holds but Student # is not a key

Course (Course#, Prof, Office, Phone)

is not in BCNF because Prof  $\rightarrow$  Office holds but Prof is not a key

But

{ (Course#, Student#, Grade), (Student#, Address, Phone) } is a BCNF decomposition of Regist.

{ (Course#, Prof), (Prof, Office, Phone) } is a BCNF decomposition of Course.

## Algorithm for BCNF decomposition

Input	U: a set of attributes F: a set of FDs
Output	D: a BCNF decomposition of U wrt F
Method	(1) $D = \{U\}$ ; (2) <b>while</b> there exists a relation scheme Q in D that is not in BCNF <b>do</b> <b>begin</b> find a nontrivial FD $X \rightarrow W$ that violates BCNF, i.e., $X \rightarrow W$ in $F^+$ and $XW \subseteq Q$ and $X \not\rightarrow Q$ ; $X^* := \{ A \mid A \text{ is in } (Q - X) \text{ and } F \models X \rightarrow A \}$ ; replace Q in D by two schemes $(X \cup X^*)$ and $(Q - X^*)$ <b>end;</b>

Note that it is NP-complete to determine whether a relation scheme is in BCNF wrt F.

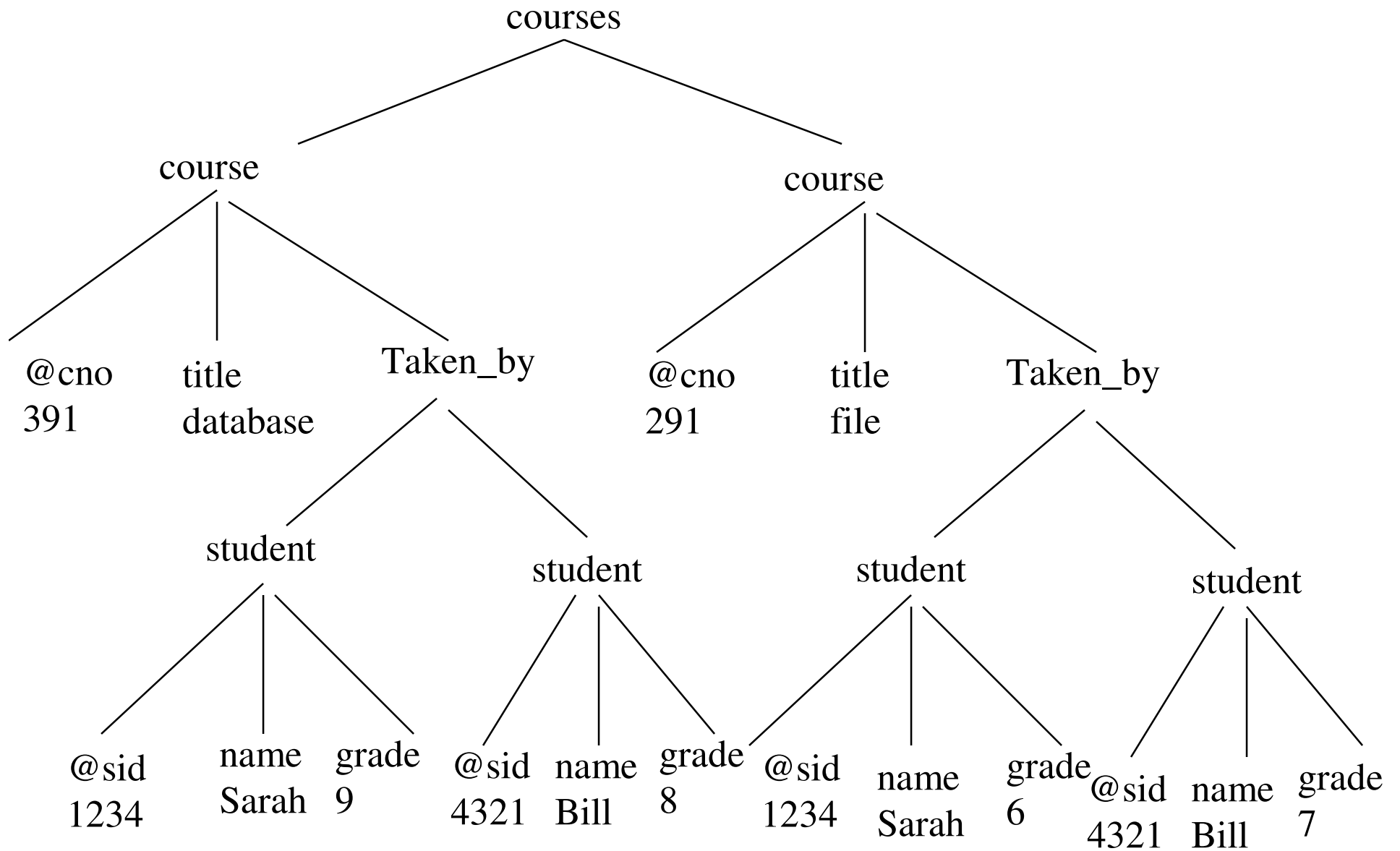


# NNF: A Normal Form for Nested Relations

- Functional dependency and multi-valued dependencies
- Path Attributes
- Minimizing redundancy and update anomalies

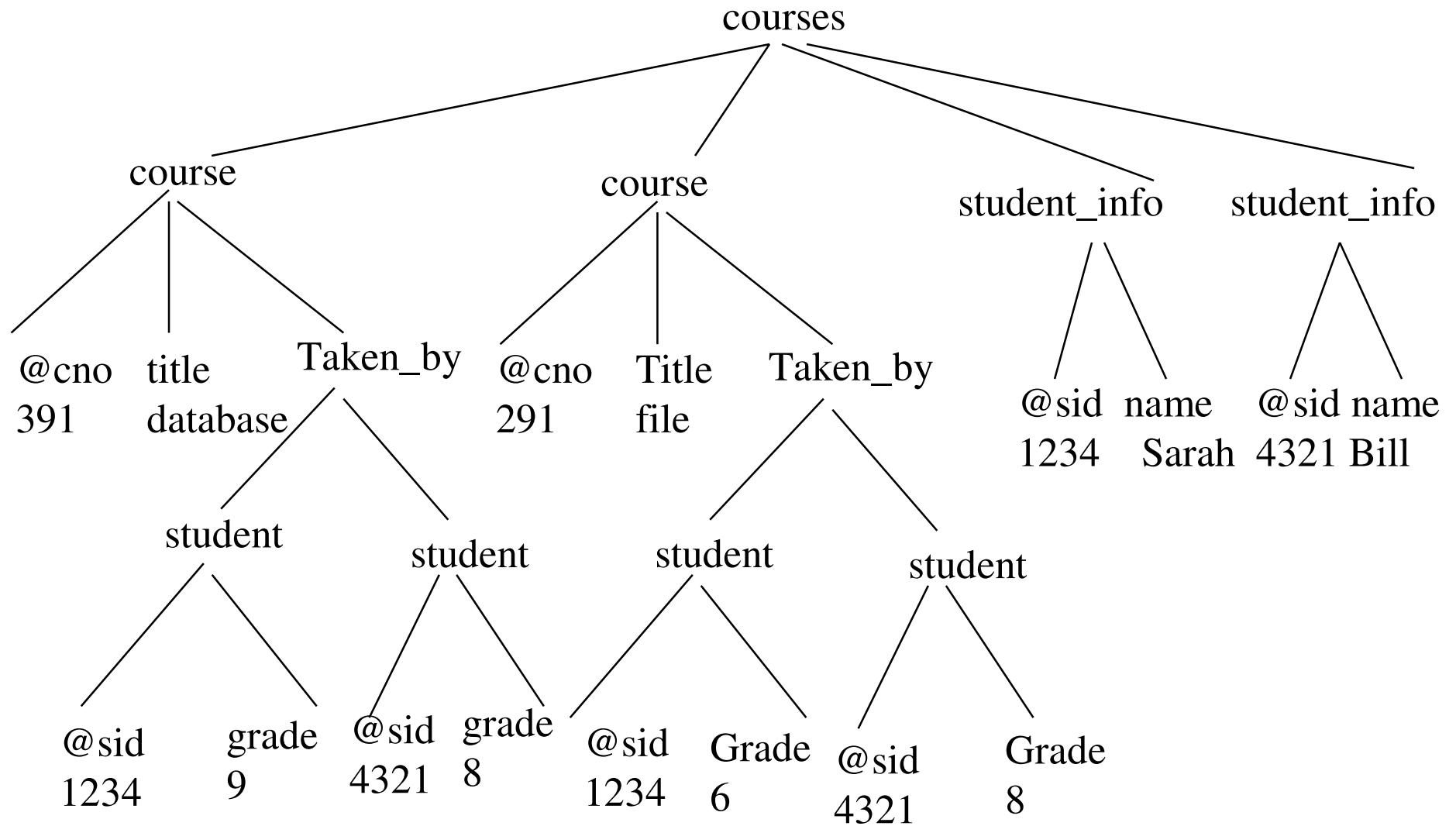
# Motivation Example for XML

```
<!DOCTYPE courses [  
  <!ELEMENT courses ( course* ) >  
  <!ELEMENT course( title, taken_by ) >  
  <!ATTLIST course cno CDATA #REQUIRED>  
  <!ELEMENT title (#PCDATA)>  
  <!ELEMENT take_by( student* )>  
  <!ELEMENT student ( name, grade)>  
  <!ATTLIST student sid CDATA #REQUIRED>  
  <!ELEMENT name ( #PCDATA)>  
  <!ELEMENT grade (#PCDATA) >  
 ] >
```



# Motivation Example for XML

```
<!DOCTYPE courses [  
  <!ELEMENT courses ( course*, student_info*) >  
  <!ELEMENT course( title, taken_by) >  
  <!ATTLIST course cno CDATA #REQUIRED>  
  <!ELEMENT title (#PCDATA)>  
  <!ELEMENT take_by( student*)>  
  <!ELEMENT student( grade) >  
  <!ELEMENT grade (#PCDATA) >  
  <!ATTLIST student sid CDATA #REQUIRED>  
  <!ELEMENT student_info( sid*, name) >  
  <!ELEMENT numberEMPTY>  
  <!ATTLIST number sid CDATA #REQUIRED>  
  <!ELEMENT name ( #PCDATA)>  
>
```



# Notations

- Assume the following disjoint sets
  - EL: the set of all element names
  - Att: the set of all attribute names, starting with @
  - Str: the set of all possible string valued attributes
  - Vert: the set of node identifies
- A DTD (Document Type Definition) is defined to be
  - $D = ( E, A, P, R, r )$ , where
    - E is a finite subset of EL
    - A is a finite subset of Att
    - P is a mapping from E to element type definitions, defined as follows
      - ◆  $P(t) = \text{EMPTY}$  or
      - ◆  $P(t) ::= \text{empty sequence} \mid t' \text{ in } E \mid P(t) \text{ union } P(t) \mid P(t) P(t) \mid P(t)^*$
    - R is a mapping from E to the power set of A
  - r is in E as the root element

- Given a DTD  $D = (E, A, P, R, t)$ , a string  $w = w_1, \dots, w_n$  is a PATH in  $D$  if
  - $w_1 = r$ ,
  - $w_i$  is in the alphabet of  $P(w_{i-1})$ , for each  $i$  in  $[2, n-1]$ , and
  - $w_n$  is in the alphabet of  $P(w_{n-1})$  or  $w_n = @l$  for some  $@l$  in  $R(w_{n-1})$
- Assume  $w$  is a path in  $D$ ,  $\text{length}(w)$  is defined as  $n$ , and  $\text{last}(w)$  as  $w_n$ .
- Given a DTD  $D$ ,
  - $\text{Paths}(D)$  stands for the set of all paths in  $D$ ,
  - $\text{Epaths}(D)$  stands for the set of all paths that ends with an element type
- DTD is recursive if  $\text{Paths}(D)$  is infinite.

# Example

```
<!DOCTYPE courses [  
  <!ELEMENT courses ( course* ) >  
  <!ELEMENT course( title, taken_by) >  
  <!ATTLIST course cno CDATA #REQUIRED>  
  <!ELEMENT title (#PCDATA)>  
  <!ELEMENT take_by( student*)>  
  <!ATTLIST student sid CDATA #REQUIRED>  
  <!ELEMENT name ( #PCDATA)>  
  <!ELEMENT grade (#PCDATA) >  
] >
```

The followings are paths in D

courses,  
courses.course  
courses.course.@cno  
courses.course.title  
courses.course.title.S  
courses.course.taken\_by  
courses.course.taken\_by.student

courses.course.taken\_by.student.@sid  
courses.course.taken\_by.student.name  
courses.course.taken\_by.student.name.S  
courses.course.taken\_by.student.grade  
courses.course.taken\_by.student.grade.S



- An XML tree  $T$  is defined to be a tree  $(V, \text{lab}, \text{ele}, \text{att}, \text{root})$ , where
  - $V$  is a finite subset of  $\text{Vert}$  ( nodes)
  - $\text{lab}: V \Rightarrow \text{EL}$
  - $\text{ele}: V \Rightarrow \text{Str} \cup V^*$
  - $\text{att}$  is a partial function  $V \times \text{Att} \Rightarrow \text{Str}$
  - $\text{root}$  in  $V$  is called the root of  $T$
- Given an XML tree  $T$ , a string  $w_1 \dots w_n$ , where with  $w_i, 1 \leq i \leq n-1$ , in  $\text{EL}$ , and  $w_n$  is in the union of  $\text{El}$ ,  $\text{Att}$ , and  $\{S\}$ .
  - The string is a path in  $T$  if there are vertices  $v_1, \dots, v_{n-1}$  in  $V$  such that
    - ⇒  $v_1 = \text{root}$ ,  $v_{i+1}$  is a child of  $v_i$  for  $1 \leq i \leq n-1$ ,  $\text{lab}(v_i) = w_i$  for  $1 \leq i \leq n-1$
    - ⇒ if  $w_n$  in  $\text{El}$ , then there is a child  $v_n$  of  $v_{n-1}$  such that  $\text{lab}(v_n) = w_n$ .
    - ⇒ If  $w_n = @l$  then  $\text{att}(v_{n-1}, @l)$  is defined
    - ⇒ if  $w_n = S$  ( $\#PCDATA$ ) then  $v_{n-1}$  has a child in  $\text{Str}$ .

- T is compatible with D if and only if
  - $\text{paths}(T)$  is a subset of  $\text{paths}(D)$

# Tree Tuples

- XML trees are defined as sets of tree tuples
- Given a DTD  $D = (E, A, P, R, r)$ , a *tree tuple*  $t$  in  $D$  is defined as a function from  $\text{paths}(D)$  to  $\text{Vert} \cup \text{Str} \cup \{\text{null}\}$  such that
  - For  $p$  in  $\text{EPaths}(D)$ ,  $t(p)$  is in  $\text{Vert} + \{\text{null}\}$ , and  $t(r) \neq \{\text{null}\}$
  - For  $p$  in  $\text{paths}(D) = \text{EPahths}(D)$ ,  $t(p)$  is in  $\text{Str} + \{\text{null}\}$ .
  - If  $t(p1) = t(p2)$  and  $t(p1)$  is in  $\text{Vert}$ , then  $p1 = p2$
  - If  $t(p1) = \text{null}$ , and  $p1$  is a prefeix of  $p1$ , then  $t(p2) = \text{null}$ .
  - $\{ p \text{ in } \text{paths}(D) \mid t(p) \neq \text{null} \}$  is finite.
- $T(D)$  is defined to be the set of all tree tuples in  $D$ .

# Example

```
<!DOCTYPE courses [  
  <!ELEMENT courses ( course*) >  
  <!ELEMENT course( title, taken_by) >  
  <!ATTLIST course cno CDATA #REQUIRED>  
  <!ELEMENT title (#PCDATA)>  
  <!ELEMENT take_by( student*)>  
  <!ATTLIST student sid CDATA #REQUIRED>  
  <!ELEMENT name ( #PCDATA)>  
  <!ELEMENT grade (#PCDATA) >  
] >
```

The followings are paths in D

t(courses) = v0

t(courses.course) = v1

t(courses.course.@cno) = 391

t(courses.course.title) = v2

t(courses.course.title.S) = database

t(courses.course.taken\_by) = v3

t(courses.course.taken\_by.student) = v4

t(courses.course.taken\_by.student.@sid) = 1234

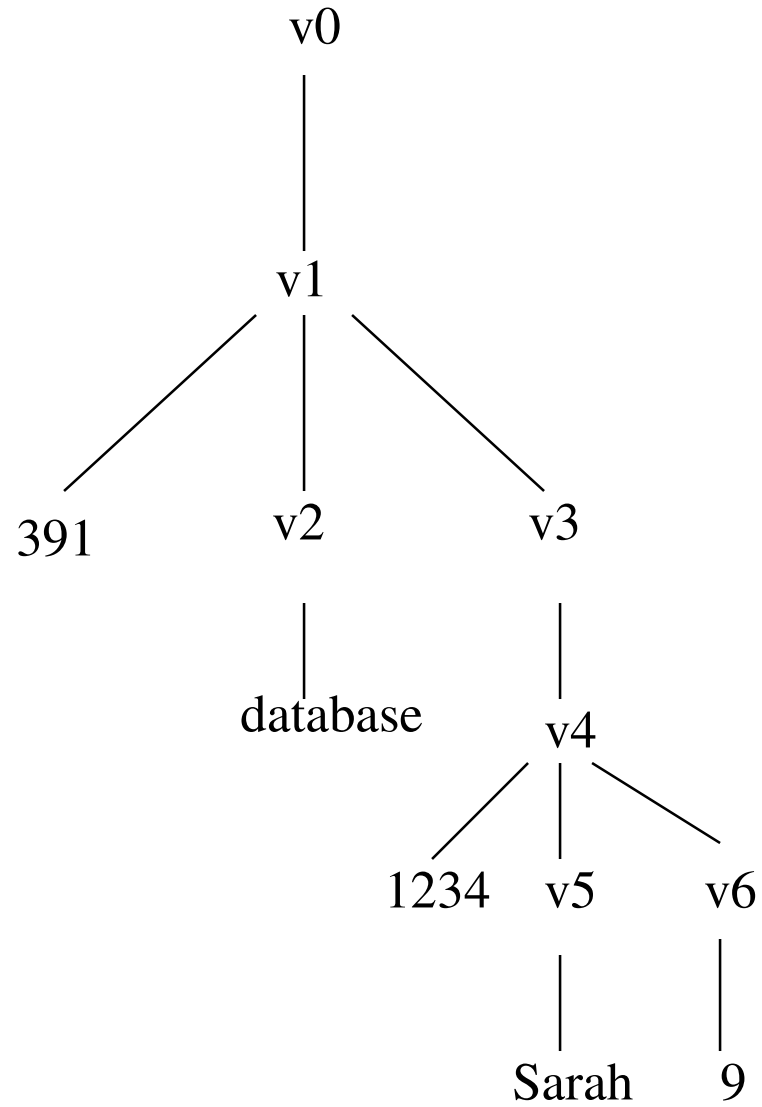
t(courses.course.taken\_by.student.name) = v5

t(courses.course.taken\_by.student.name.S) = Sarah

t(courses.course.taken\_by.student.grade) = v6

t(courses.course.taken\_by.student.grade.S) = 9

# The XML tree for this one tree tuple



## ■ Important Results:

- Given a DTD  $D$  and an XML tree  $T$  such that  $T$  conforms with  $D$ . Then  $T$  can be represented by a set of tree tuples, if we consider it as an unordered tree.

# Functional Dependencies

- Let  $D$  be a DTD,  $S1$  and  $S2$  are finite non-empty subsets of  $\text{paths}(D)$ .
  - A functional dependency FD over  $D$  is an expression of the form  **$S1 \twoheadrightarrow S2$**
  - An XML tree  $T$  satisfies  $S1 \twoheadrightarrow S2$  if for every pair of tree tuples  $t1, t2$  in  $\text{tuples}(T)$ ,
    - $t1.S1 = t2.S2$  and  $t.S1 \neq \text{null}$  implies  $t1.S2 = t2.S2$ .

# Example

```
<!DOCTYPE courses [  
  <!ELEMENT courses ( course* ) >  
  <!ELEMENT course( title, taken_by ) >  
  <!ATTLIST course cno CDATA #REQUIRED>  
  <!ELEMENT title (#PCDATA)>  
  <!ELEMENT take_by( student* )>  
  <!ATTLIST student sid CDATA #REQUIRED>  
  <!ELEMENT name ( #PCDATA)>  
  <!ELEMENT grade (#PCDATA) >  
 ] >
```

The followings are paths in D

courses,

courses.course

courses.course.@cno

courses.course.title

courses.course.title.S

courses.course.taken\_by

courses.course.taken\_by.student

courses.course.taken\_by.student.@sid

courses.course.taken\_by.student.name

courses.course.taken\_by.student.name.S

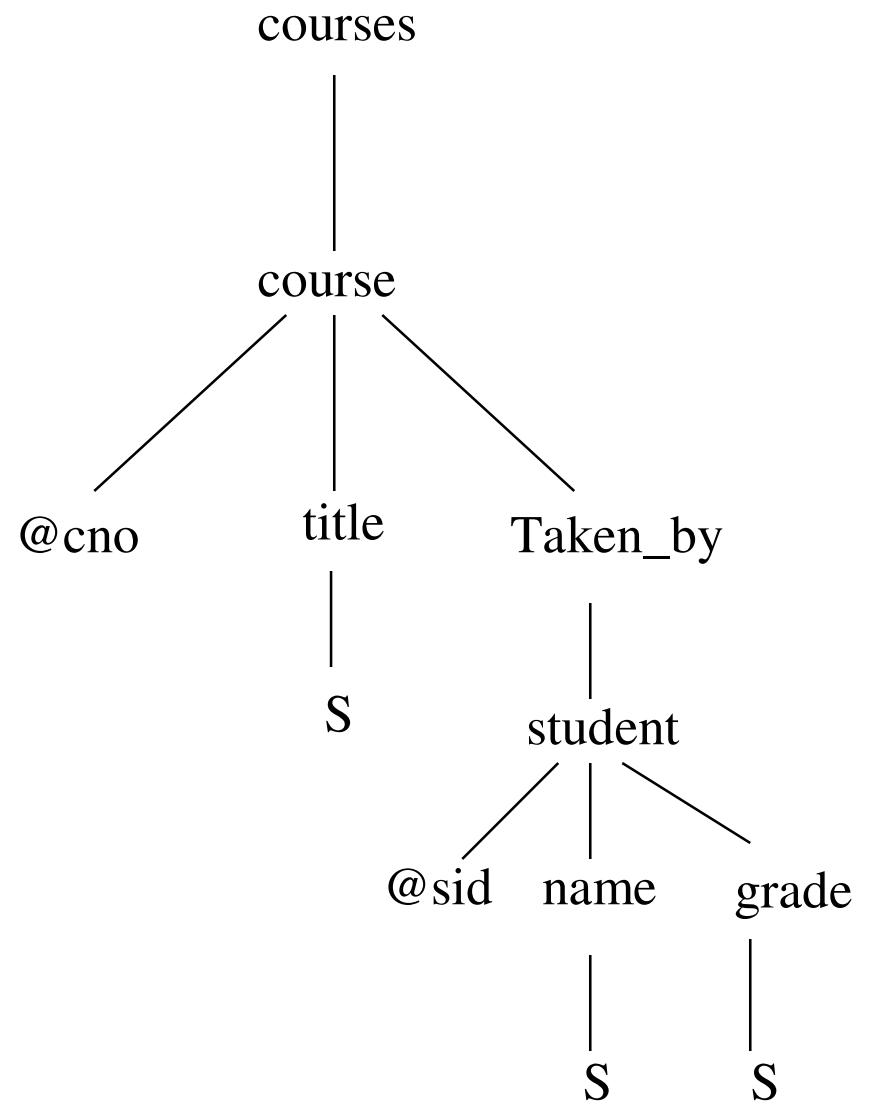
courses.course.taken\_by.student.grade

courses.course.taken\_by.student.grade.S

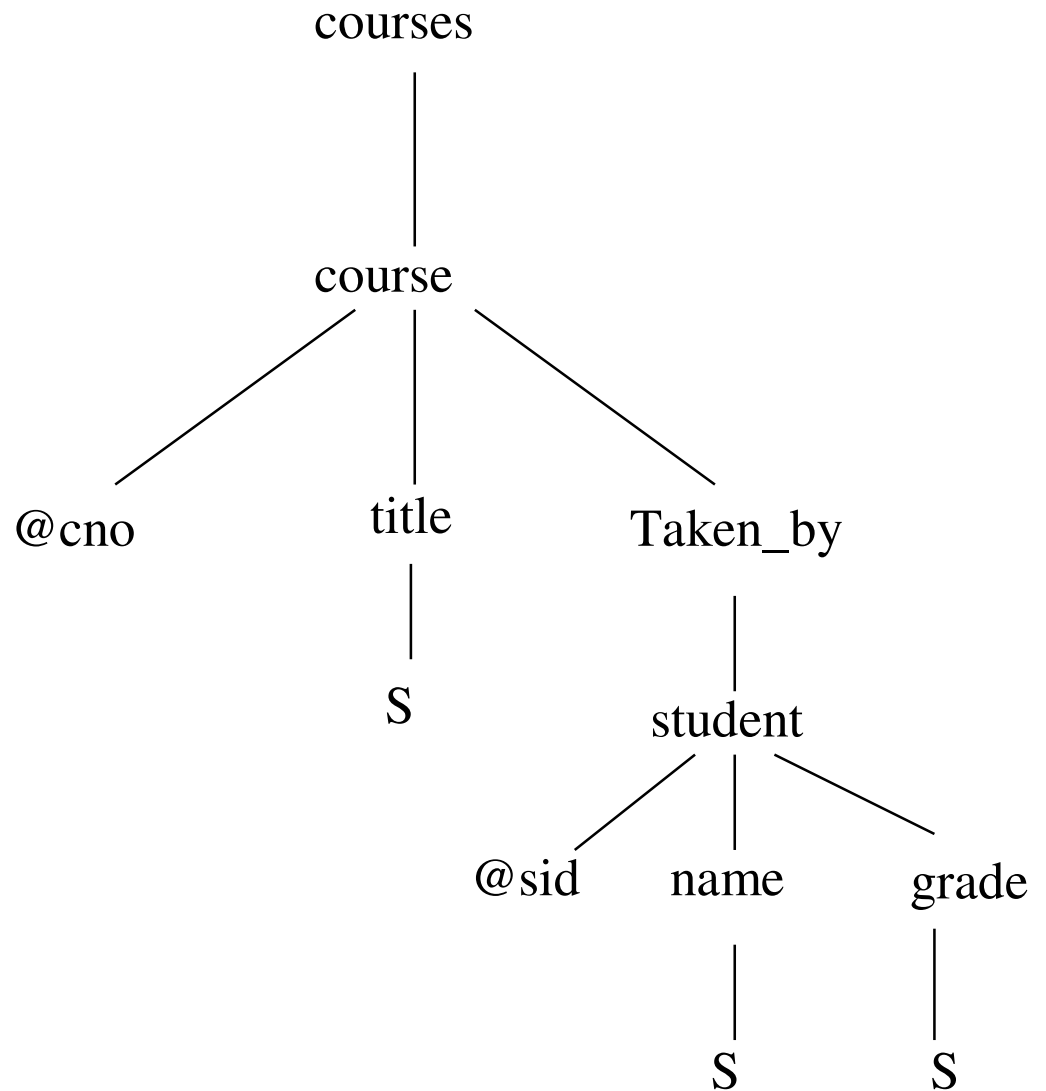


# Example: Paths(D)

courses,  
courses.course  
courses.course.@cno  
courses.course.title  
courses.course.title.S  
courses.course.taken\_by  
courses.course.taken\_by.student  
courses.course.taken\_by.student.@sid  
courses.course.taken\_by.student.name  
courses.course.taken\_by.student.name.S  
courses.course.taken\_by.student.grade  
courses.course.taken\_by.student.grade.S



# Example



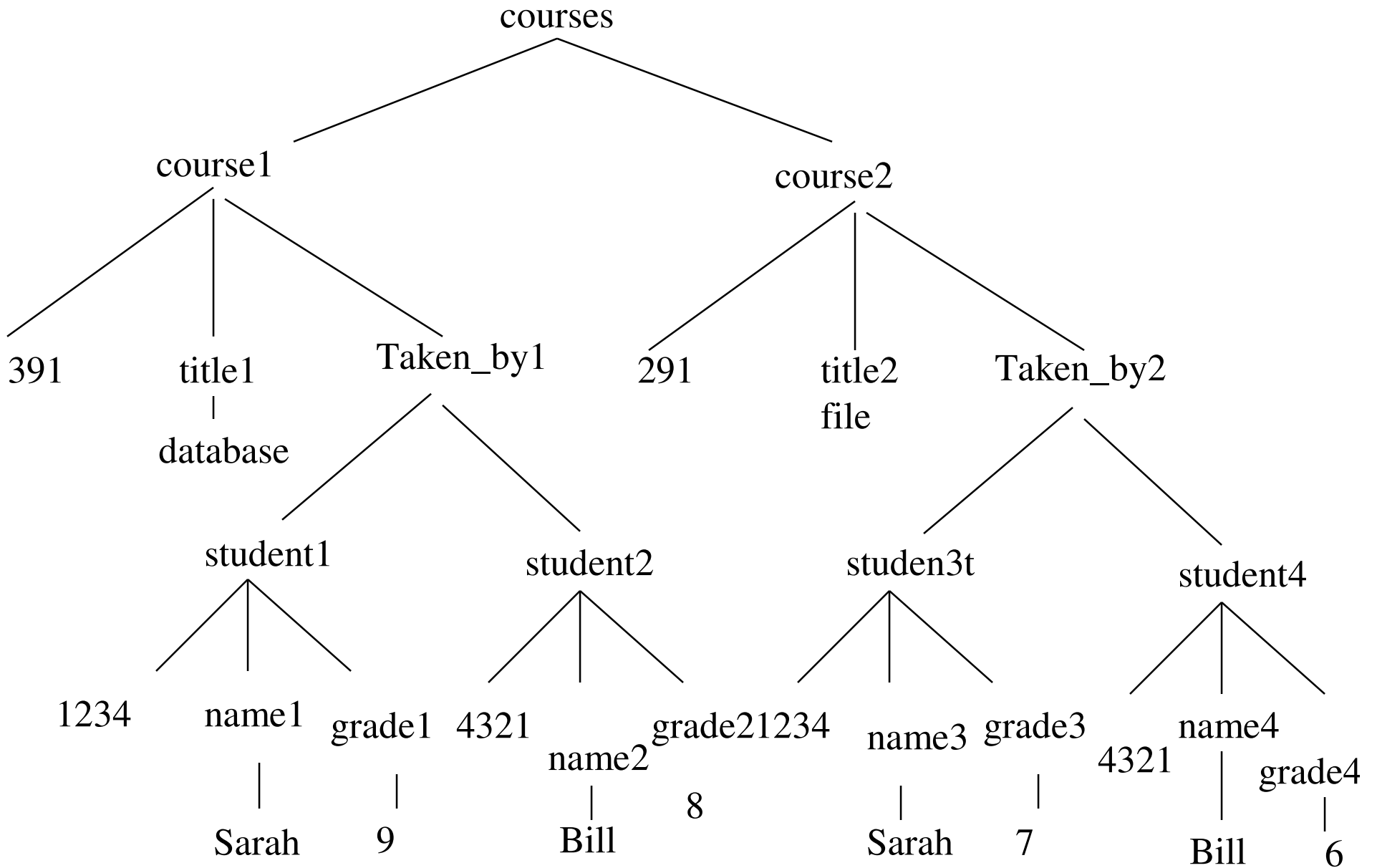
## Constraint:

cno is a key of course

## FD1:

**courses.course.@cno --> courses.course**

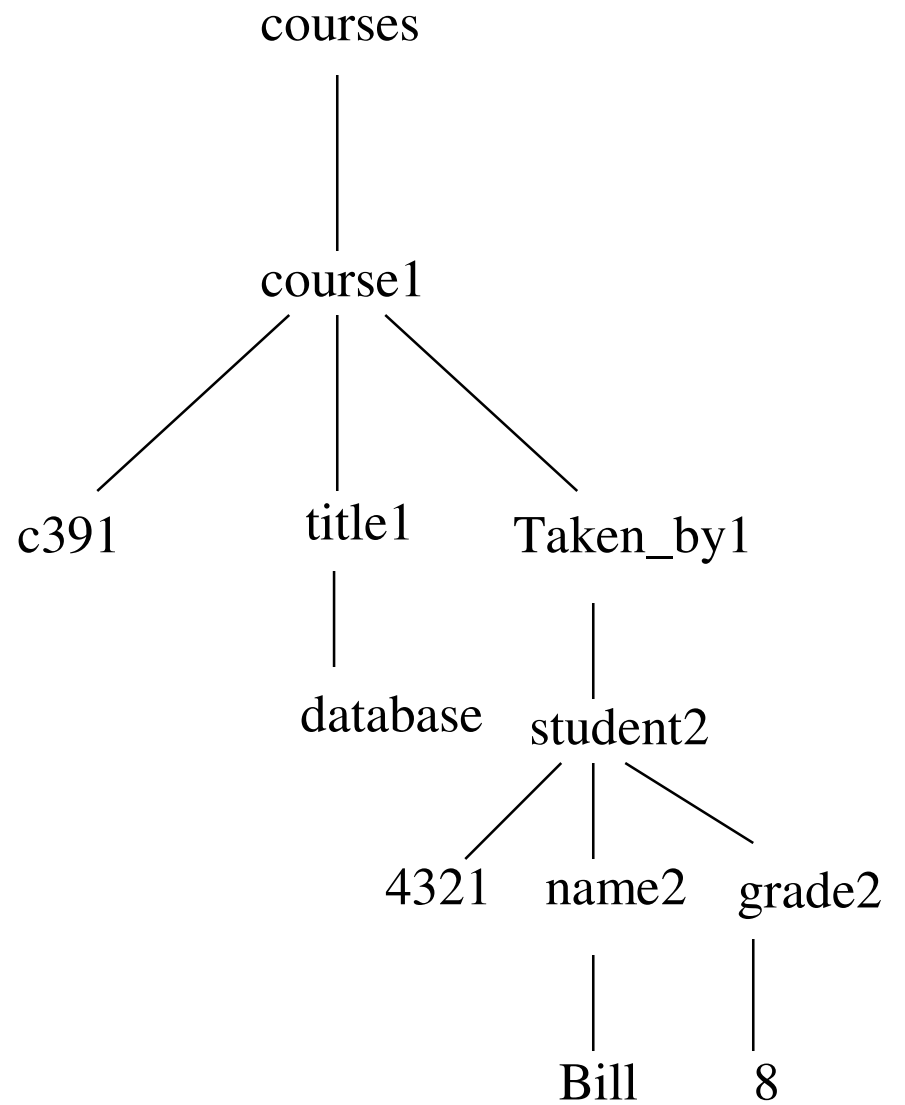
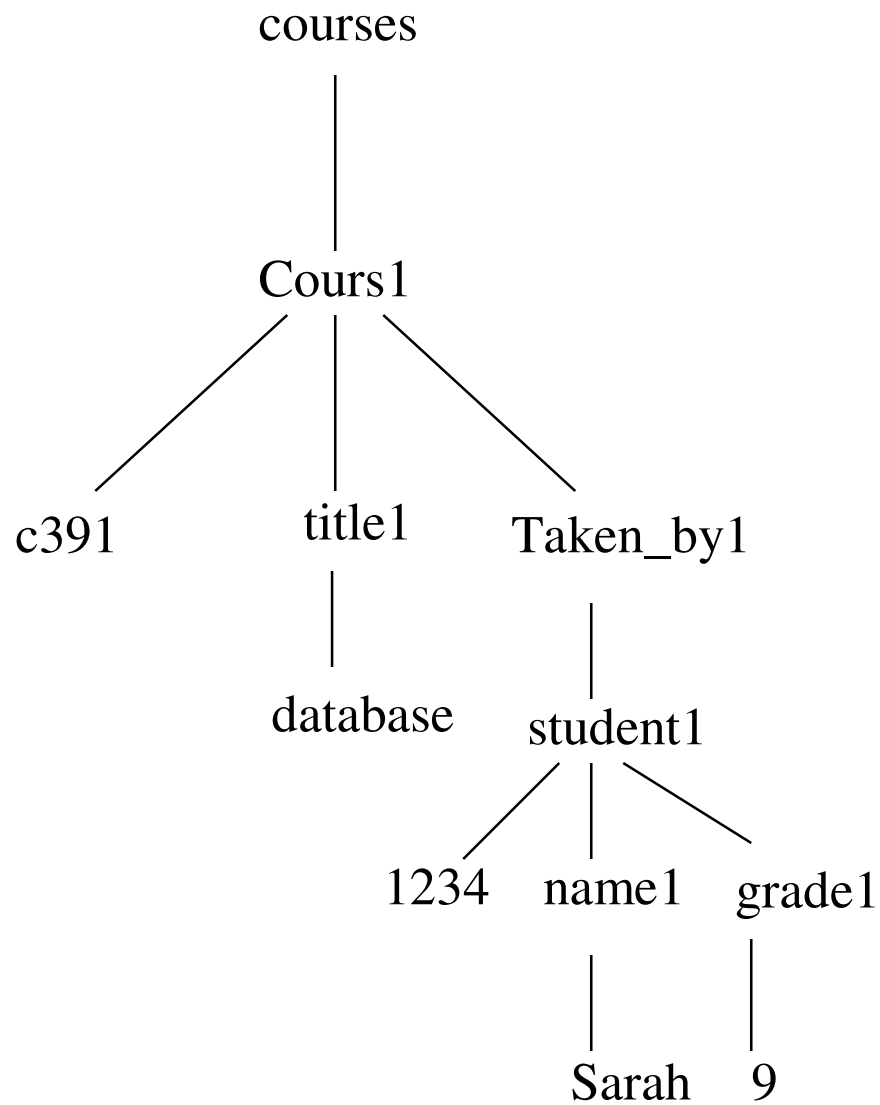
# Consider a sample XML Tree T1



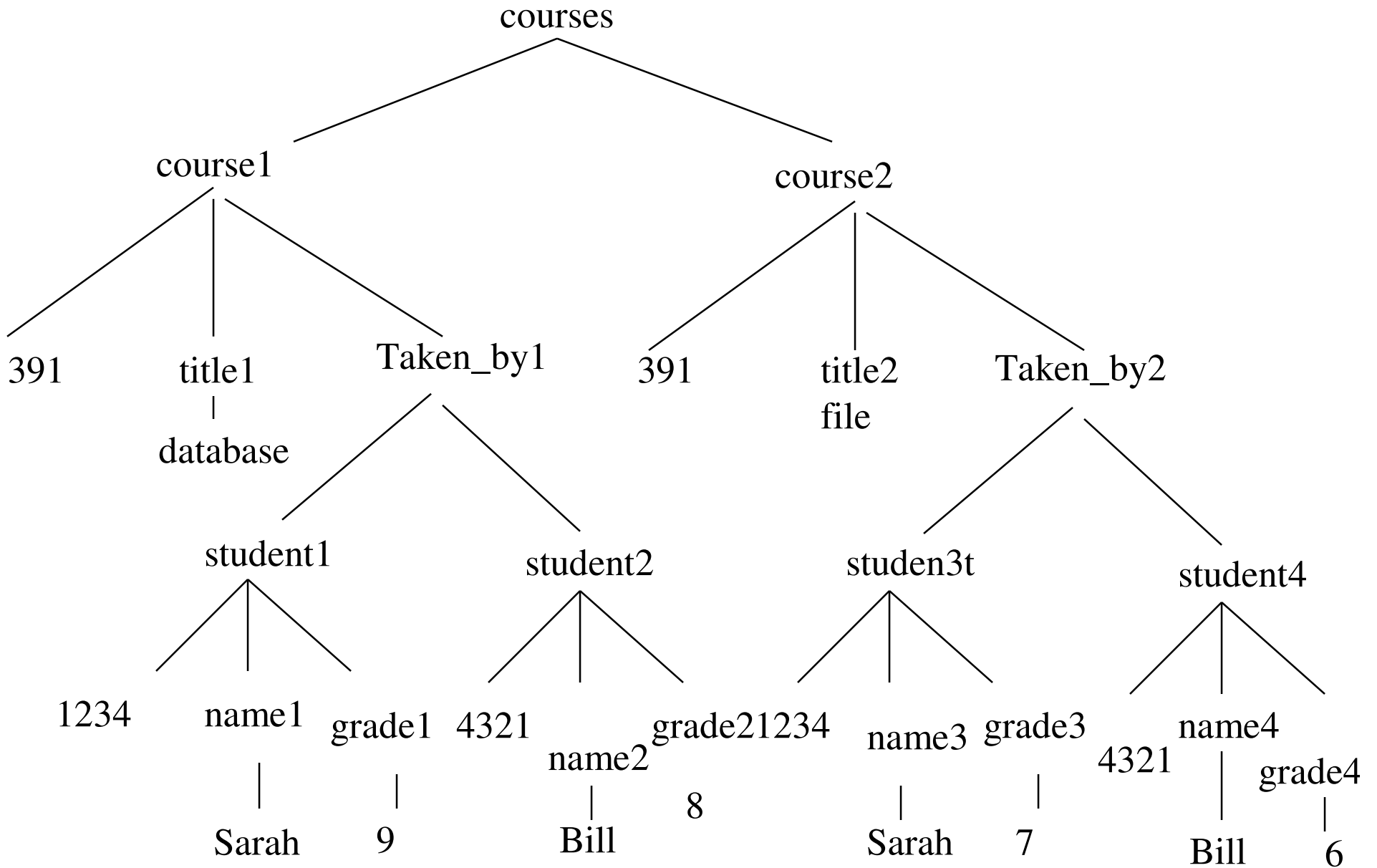
## The corresponding flat table for T1

<b>Cno</b>	<b>Title</b>	<b>SID</b>	<b>Name</b>	<b>Grade</b>
391	database	1234	Sarah	9
391	database	4321	Bill	8
291	file	1234	Sarah	7
291	file	1234	Bill	6

The following are the only two tree types with cno = c391 in T1



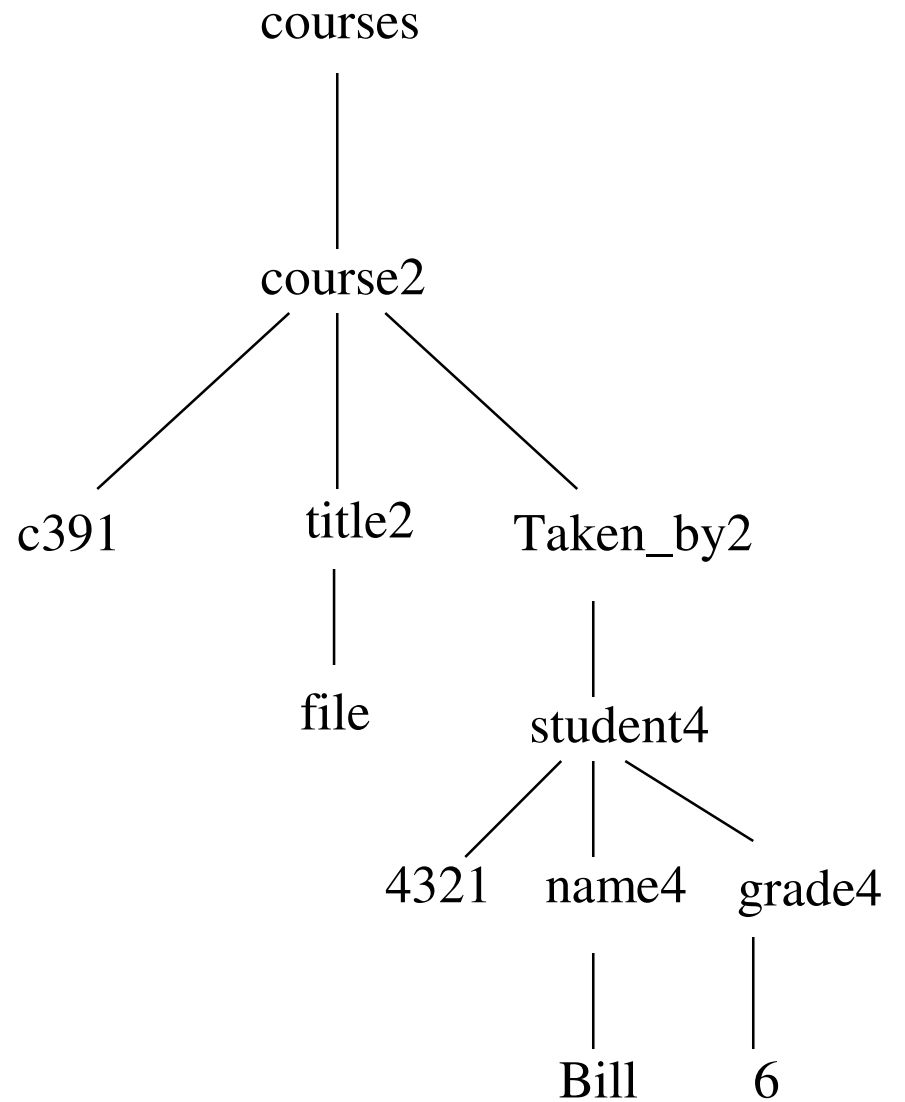
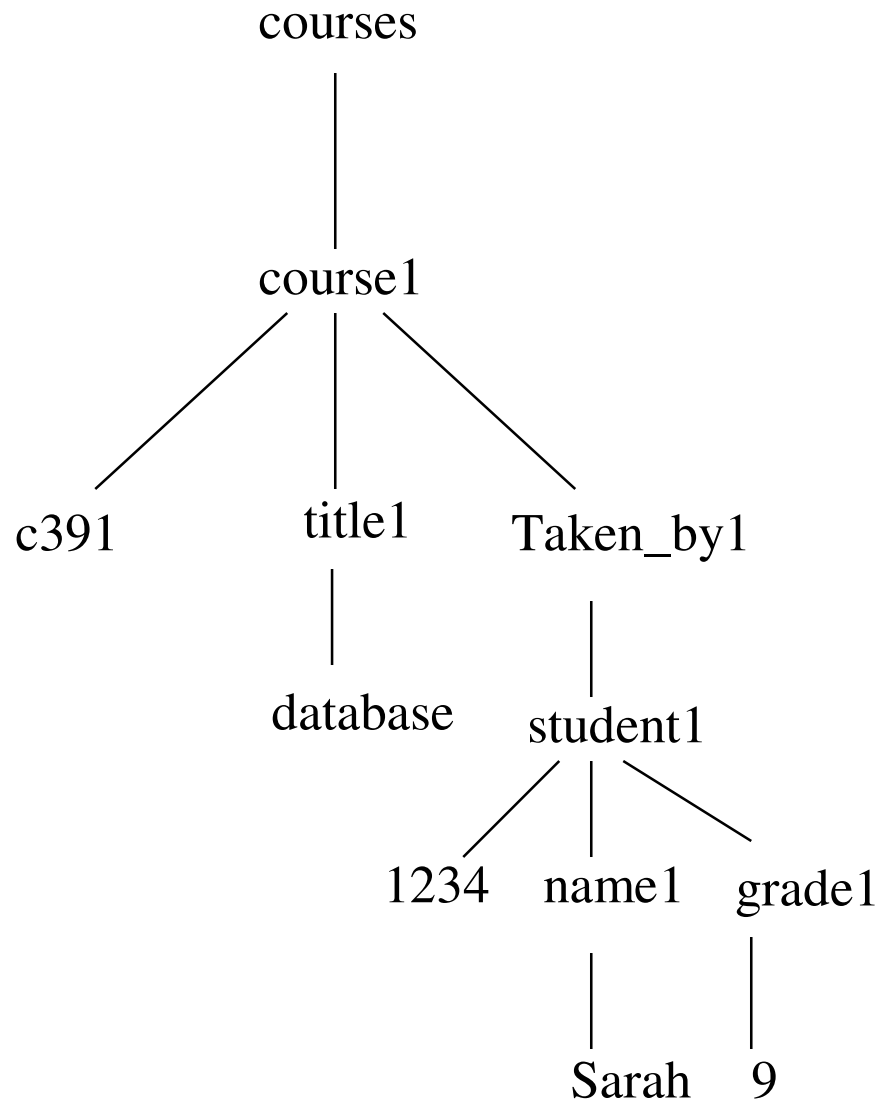
# Consider another XML Tree T2



The corresponding flat table for T2

<b>Cno</b>	<b>Title</b>	<b>SID</b>	<b>Name</b>	<b>Grade</b>
391	database	1234	Sarah	9
391	database	4321	Bill	8
391	file	1234	Sarah	7
391	file	1234	Bill	6

The following are two tree types with  $cno = c391$  in T2





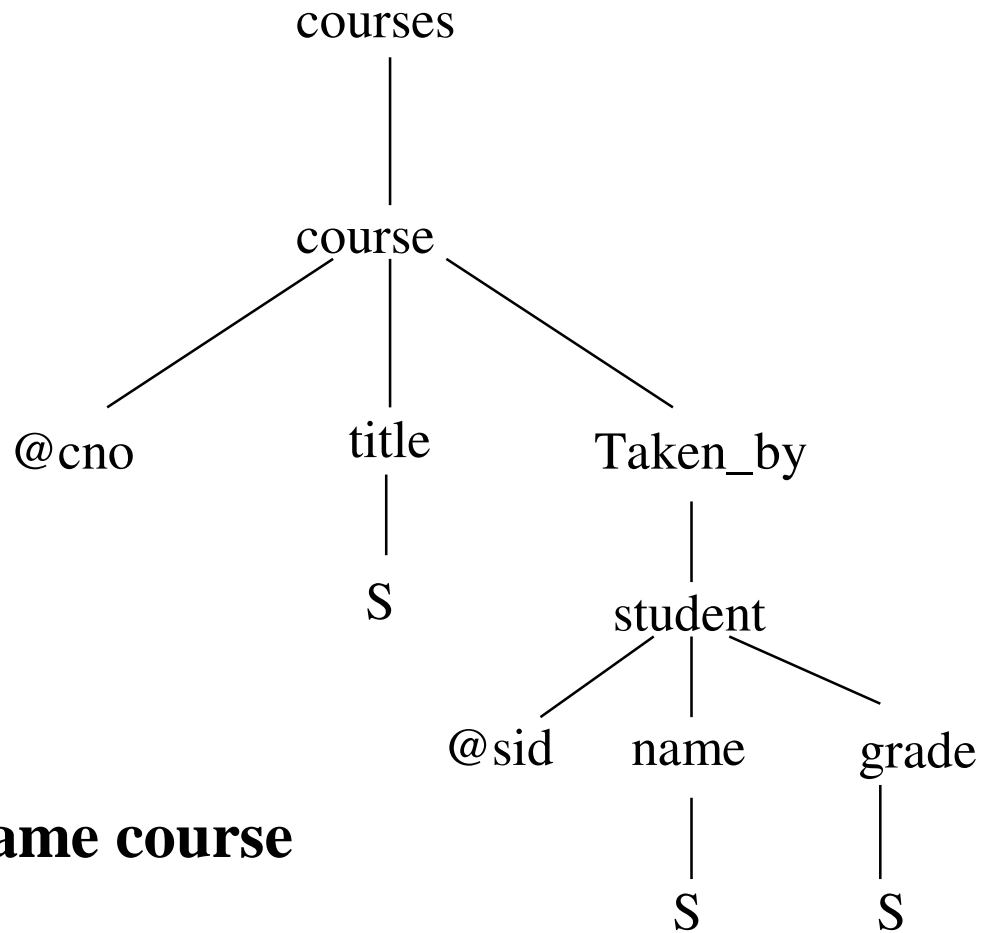
## ■ Observation

- Both T1 and T2 conform to the DTD
- T1 satisfies the FD

**|||→courses.course.@cno --> courses.course**

- T2 does not satisfy the above FD

# Example



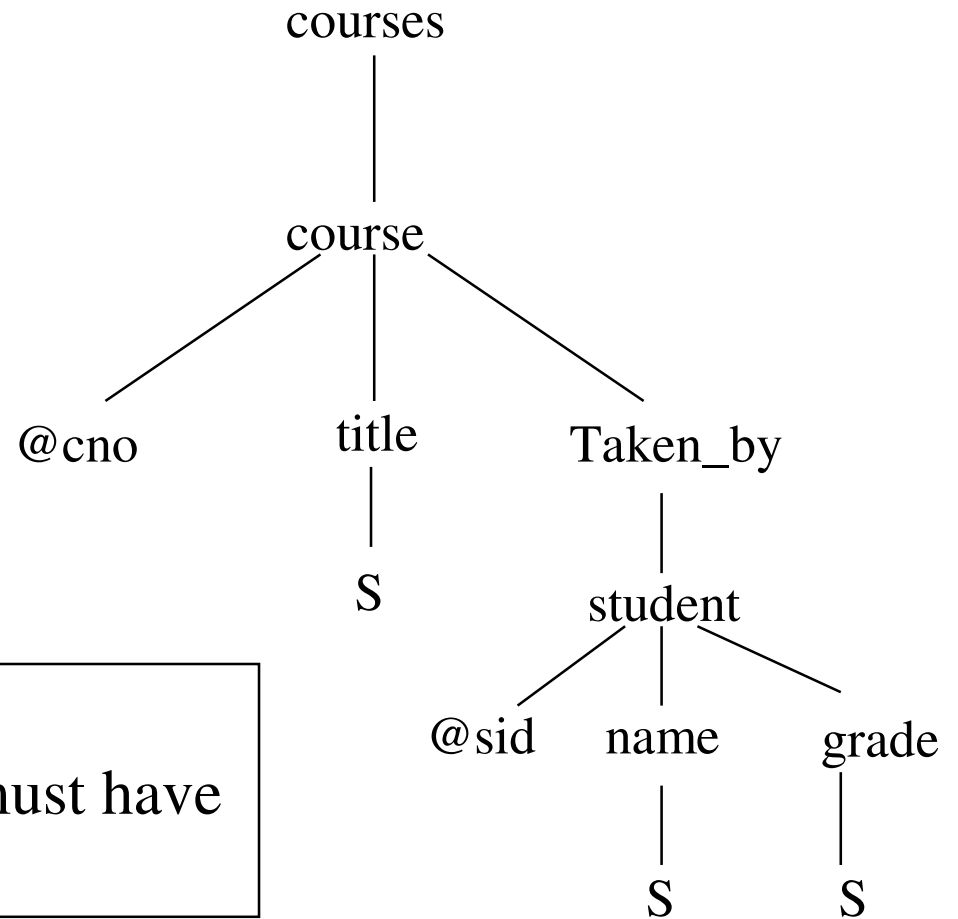
## Constraint:

two distinct students of the same course cannot have the same sid

## FD2:

{ courses.course, courses.course.taken\_by.student.@sid }  
--> courses.course.taken\_by.student

# Example



**Constraint:**

two students with the same sid must have the same name

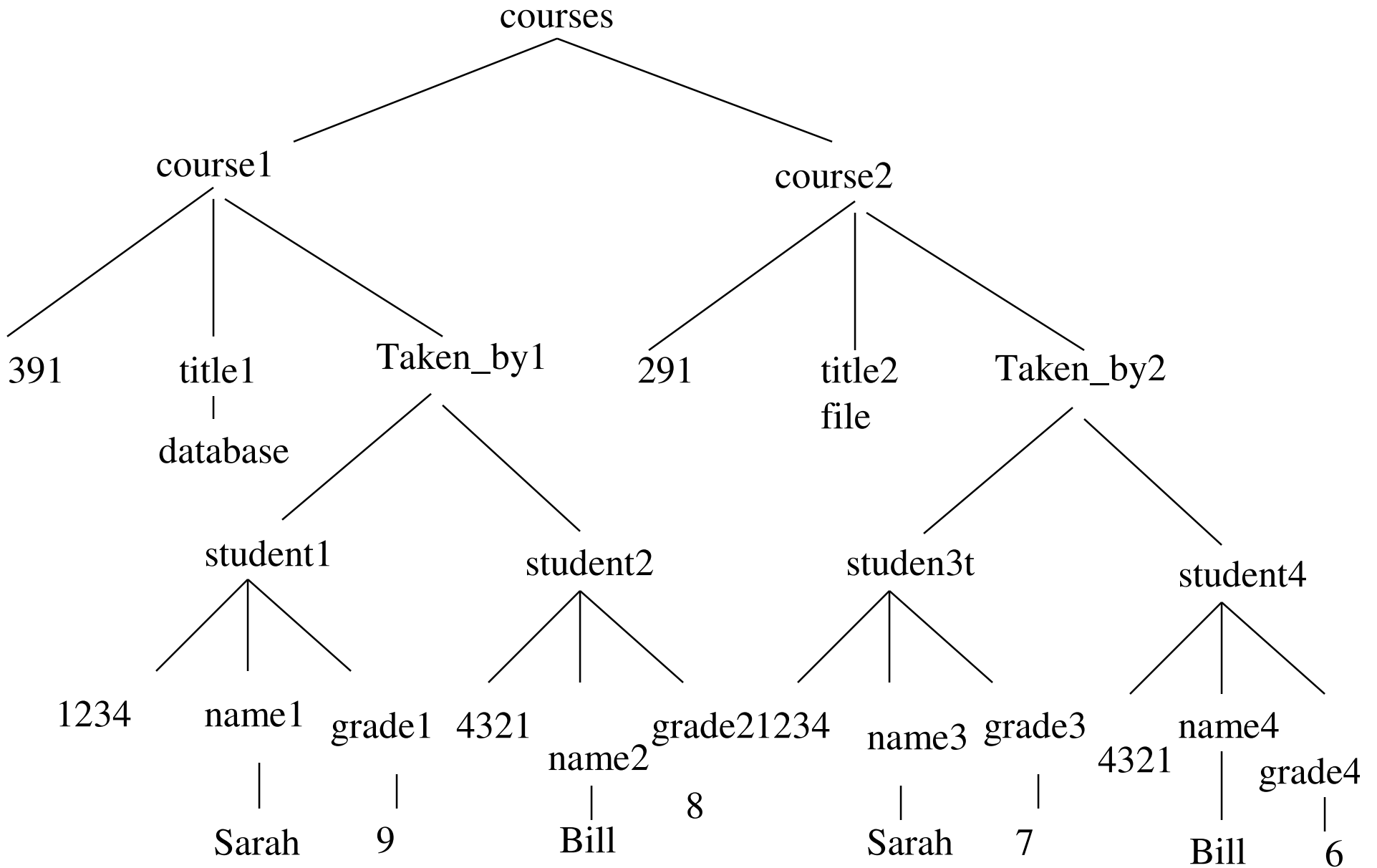
**FD3:**

`courses.course.taken_by.student.@sid -->`  
`courses.course.taken_by.student.name.S`

# XNF: An XML Normal Form

- Given a DTD, and a set  $F$  of FDs,  $(D, F)$  is in XML normal form (XNF) if and only if for every nontrivial FD of the form  $S \twoheadrightarrow p.@l$  or  $S \twoheadrightarrow p.S$ , it is the case that  $S \twoheadrightarrow p$  is implied by  $F$ .
- Intuition
  - For every set values of the elements in  $S$ , we can find only one value of  $p.@l$ . Thus, we need to store the value only one.

Consider the following example again



**We have FD3:**

```
courses.course.taken_by.student.@sid -->  
courses.course.taken_by.student.name.S
```

But the following does not hold:

```
courses.course.taken_by.student.@sid -->  
courses.course.taken_by.student.name
```

This implies that the student name for a given sid, the document may have multiple copies of student name.

# Relationships with other normal forms

- Assume a standard coding between tables and XML documents
  - A relation schema is in BCNF if and only if its XML counterpart is in XNF
- Assume a standard nesting operations and coding
  - A nested relation is in NNF if and only if its XML representation is in XNF.

# Normalization Algorithm

- Two basic operations
  - Moving attributes
  - Creating new element types
- Given a DTD  $D$  and a set  $F$  of FDs
  - If  $(D, F)$  is in XNF, return
  - Otherwise find an anomalous FD and use the two basic operations to modify  $D$  to eliminate the anomalous FD,
  - Continue the above steps until  $(D, F)$  is in XNF.
- The normalization algorithm is efficient and join-lossless