

CO130 - Databases

Prelude

The content discussed here is part of CO130 - Databases (Computing MEng); taught by Thomas Heinis, in Imperial College London during the academic year 2018/19. The notes are written for my personal use, and have no guarantee of being correct (although I hope it is, for my own sake).

Material Order

These notes are primarily based off the slides on CATE. This is the order in which they are uploaded (and I'd assume the order in which they are taught).

1. *Introduction.pdf*
2. *ER Modelling.pdf*
3. *ER to RM.pdf*
4. *Relational Algebra.pdf*
5. *Tutorial ER.pdf*
6. *Tutorial Translation.pdf*
7. *Relational Algebra.pdf*
8. *Functional Dependencies.pdf*
9. *Tutorial Relational Algebra.pdf*
10. *Normalisation.pdf*
11. *SQL.pdf*
12. *Data Definition.pdf*
13. *Data Manipulation.pdf*
14. *Advanced SQL.pdf*
15. *Functional Dependency Tutorial.pdf*
16. *Transactions.pdf*
17. *SQL Tutorial.pdf*
18. *Storage.pdf*
19. *Indexing.pdf*
20. *NoSQL.pdf*
21. *MapReduce.pdf*

Introduction

We use databases as it's more organised; hence it's easier to model and manage. It's more efficient, as it's fast to search, and update, and integration allows us to minimise data duplication. Concurrent (and therefore multi-user) access allows multiple people to access the database at the same time (will require some techniques).

Transactions are sequences of database actions that execute in a coherent, and reliable way - the classical properties are **ACID**. Consider the two transactions T1: $A = A - 100; B = B + 100$, and T2: $B = B - 100; A = A + 100$, we can observe ACID properties as follows;

- **atomicity** if one part of a transaction fails, the entire transaction fails on completion of T1, either $A' = A - 100$, and $B' = B + 100$, or $A' = A$, and $B' = B$, where the former is a successful transaction, and the latter is in the case of a failure.

- **consistency** transactions don't leave the database in an inconsistent state
the sum of the balances must remain the same, such that $A' + B' = A + B$; we can also have more constraints such as keeping balances positive, or limiting the amount a transfer can do at once
- **isolation** transactions run as if no other transactions are running (may need to wait)
given the two concurrent transactions T1, and T2, one has to be completed before the other can start
- **durability** results of successful transactions aren't lost on system failure
the new values, A' , and B' must persist if the transaction completes, even if the system fails (disk failure etc.)

A **Database Management System (DMBS)** creates new databases via a **Data Definition Language (DDL)**, which specifies the structure (**schema**). It also queries, and manipulates through a **Data Manipulation Language (DML)**. Examples of this include *PostgreSQL*, *MySQL*, *SQLite*, and can also fall under *NoSQL*, however SQL remains as the most widespread technology (as of writing this).

It lets us define, query, and manipulate databases with a high-level declarative language (**Structured Query Language**). It's standardised by the ISO, but each DBMS implements its own variation of the standards, which may be costly if it's complex.

Relational Model

Consider the following model, represented as a table;

heading	title:string	year:int	length:int	genre:string
body	Gone with the Wind	1939	231	Drama
	Star Wars	1977	124	Science Fiction
	Wayne's World	1992	95	Comedy

We have the columns be the attribute, with the top row being the heading, and the rest being the body. The attributes are in the format **name:type**. The rows (of the body) are referred to as tuples. A relation is the heading as well as the body (the entire table). The heading is an **unordered set** of attributes, and an attribute is the name as well as the type (typically indivisible types). The body is an unordered set of tuples, and a tuple is the set of attribute values. The schema is for the entire relation is the name of the relation, and the heading, in this case, we'd have; **movies(title:string, year:int, length:int, genre:string)**, and a database is a collection of relations. A schema for a database is the schemas for all relations.

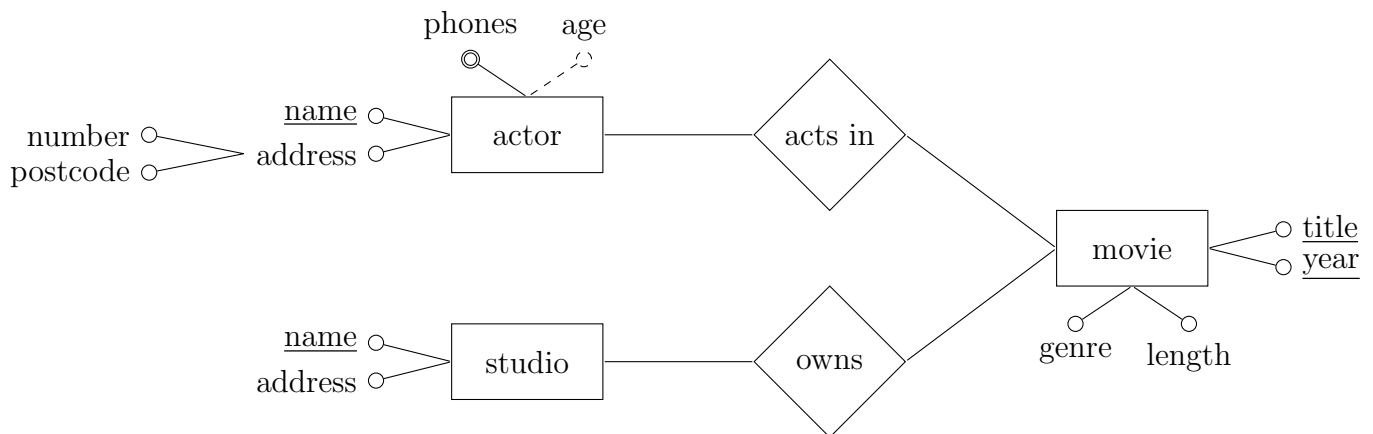
In mathematics, with a set of sets; S_1, S_2, \dots, S_n , a relation R is a set of tuples T_1, T_2, \dots, T_n , where $T_k \in S_k$, therefore $R \subset S_1 \times S_2 \times \dots \times S_n$. As the idea of relations stems from set theory, it's important to note that the order in which we represent the attributes, and tuples in unimportant. In the Relational model we have **attributed** tuples, rather than **ordered**; R is the set of tuples $(A_1 : S_1 = T_1, \dots, A_n : S_n = T_n)$, with $T_k \in S_k$. It's important to note that relations aren't 2-dimensional tables, even though it's more convenient to draw it on paper. We should instead consider them as a set of n -dimensional values, such that we have (**title:string=StarWars, year:int=1997, length:int=127, genre:string=ScienceFiction**), as a 4-dimensional movie value.

Entity Relationship Modelling

When a new database is being developed, it's important to try and model the real-world situation, instead of trying to refine it into an implementation, such as a relational model. In Entity-Relationship modelling, we try to create a diagram which represents the information needed for the database (the Entity Relationship Diagram). As there is no universally accepting notation for ER diagrams, we will use the following notation;

type	description	shape
entity sets	a set of distinguishable entries that share the same set of properties, can be physical (a room etc.), an event (flight, sale, etc.) - they normally correspond to nouns	rectangle
relationship sets	captures how two or more entity sets are related (e.g. owns, tutors), we can also have more than one relationship set between entity sets, and they can also have a relationship set on the same entity - they sometimes correspond to verbs	diamonds
attributes	properties of an entity; relationship sets can also have attributes, and primary keys are underlined	small circles

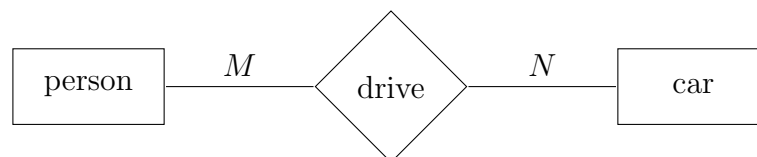
The movie example is represented below. Note that I've also extended it to contain different types of complex attributes. You can see that the address field is subdivided into number, and postcode, we have a multi-valued attribute in phones, and a derived attribute in age (can be calculated from date of birth)



Cardinality Constraints

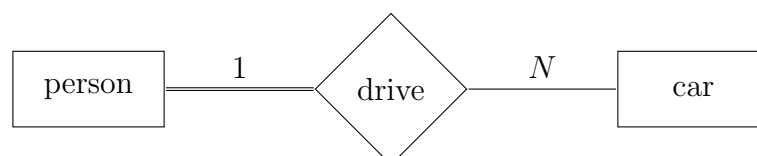
A relationship between two entity sets can be seen as one of the following (using the car example, where a person can drive N cars, and a car can be driven by M people);

- one-to-one $M = 1, N = 1$
- one-to-many $M = 1$
- many-to-many no restrictions

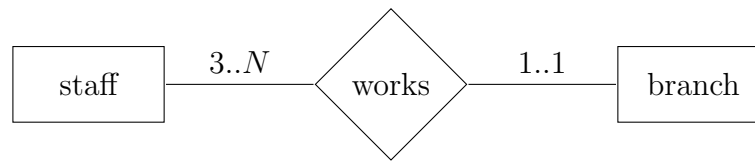


Participation

Using the same example, we can force participation, which means that all entities must participate in a relationship, with a double line. The example below, with the cars, suggests that everyone drives at least 1 car, and car can have at most 1 driver (note that it doesn't mean that every car has a driver).



Sometimes, instead of using double lines, E-R notation may allow explicit bounds. The below diagram suggests that staff work in exactly 1 branch, and that each branch must have at least 3 members of staff;



Fan, and Chasm Traps

See the diagrams in *ER Modelling.pdf* for examples.

We need to ensure that we do not allow ambiguous paths between entities. For example, if we say **staff** works for **faculty** which operates **department**, if some member of staff works for Engineering, which operates both Computing, and EE, we can't follow a path from the staff member to their department. This can be solved by having the staff work for the department, which is operated by some faculty.

If we can't follow a path between two entities, we may need to add another relation between them to specify the relationship.

Weak Entities

Entities which cannot be uniquely identified by their own attributes are called **weak** entities, in contrast to **strong** entities which have primary keys. A weak entity has to be defined by a strong entity. Weak entities are drawn as a double rectangle, and the relationship to the strong entity is drawn as a double diamond. The key for a weak entity is formed by combining the primary key of the strong entity with attributes of the weak entity (denoted by a dashed underline). For example, a room in a building is a **weak** entity, with a room number, and can be uniquely identified with the name of the building combined with the room number.

ER to RM

Keys are used to join information when we do queries. The primary key is the unique identifier of a tuple, and a foreign key is the primary key of **another** table.

Once we have an Entity Relationship Model, we can then map it to a relational schema. The schemas can then be refined with functional dependencies, and implemented with relational languages. A high level data model isn't the only concern for database design (considered a simpler aspect).

If we consider a strong entity set, with simple attributes, we can easily create a table from it (see the example, **movie**, drawn above).

```
1 movie(title, year, length, genre)
2
3 create table movie {
4     title varchar(120),
5     year int,
6     length int,
7     genre char(20),
8
9     primary key(title, year)
10 }
```

Composite attributes are also fairly easy to represent, since we just flatten it, therefore store each of the sub-attributes as their own field - consider the actor example (note that if the primary key is composite, we mark all the sub-attributes as primary keys). This example also shows how multi-valued attributes are stored - we create their own relation, and link back to the entity set with a foreign key constraint;

```

1 actor(name, number, postcode)
2 actor_phones(actorName, phoneID, other attributes)
3
4 create table actor {
5     name varchar(60),
6     number int,
7     postcode varchar(10),
8
9     primary key (name)
10 }
11
12 create table actor_phones {
13     actorName varchar(60),
14     phoneID varchar(10),
15     ...
16     primary key (actorName, phoneID),
17     foreign key (actorName) references actor.name
18 }

```

The relational model doesn't allow us to specify derived attributes, and we're therefore expected to calculate them with queries.

We can also consider how many-to-many relationships can be represented (using the car example again). The first example is when a car can be driven by many people, and a person can drive many cars;

```

1 person(ID, other attributes)
2 car(regno, other attributes)
3 drive(personID, regno, other attributes)
4
5 create table drive {
6     personID varchar(10),
7     regno varchar(12),
8     ...
9     primary key (personID, regno),
10    foreign key (personID) references person.ID on delete cascade,
11    foreign key (regno) references car.regno on delete cascade
12 }

```

However, it's easier to represent a one-to-many relationship, as we can simply include the primary key of the "one" relation as a foreign key in the many (where a car can be driven by 1 person, and a person can drive many cars). The representation for a one-to-one is similar, but the foreign key can be in either one, which is up to the designer;

```

1 person(ID, other attributes)
2 car(regno, personID, other attributes)
3
4 create table car {
5     regno varchar(12),
6     personID varchar(10),
7     ...
8     primary key (regno),
9     foreign key (personID) references person.ID
10 }

```

In a weak entity, we include the primary key of the strong relation as a foreign key, but also use it as part of the composite primary key for the weak entity. A delete cascade is also used, and finally the we use `not null` to ensure total participation;

```

1 building(name, other attributes)
2 room(no, buildingname, other attributes)
3
4 create table room {
5     no varchar(120),
6     buildingname varchar(50) not null,
7     ...
8     primary key (no, buildingname),
9     foreign key (buildingname) references building.name on delete cascade
10 }

```

A multiway relationship can be mapped as several binary relationships, or we can generalise it to have the primary key consist of the relationship be composed of the primary keys of the many entity sets, and have the primary keys of the one entities be attributes. For roles, we map each role as a foreign key attribute in the entity set;

```

1 movie(ID, other attributes)
2 sequelof(originalID, sequelID, other attributes)
3
4 create table sequelof {
5     originalID int,
6     sequelID int,
7     ...
8     primary key (originalID, sequelID),
9     foreign key originalID references movie.ID,
10    foreign key sequelID references movie.ID
11 }

```

Extended Models

Newer models extend the classic model by supporting features from **object-oriented design**. We can look at how E-R models handle specialisation, or generalisation with **is-a hierarchies**. For example, we can specify that a **cartoon is-a movie**, which means that it inherits all the attributes of movie, but can have more attributes. and relations. Once this extends on multiple levels, we begin to form a hierarchy.

Relational Algebra

We can use **relational algebra** to construct new relations from existing ones, the operators include selection, projection, intersection, union, difference, and product. Consider the examples below, **a**, **b**, and **c** are all integers. Intersection, union, and difference will be omitted, since those are fairly self-explanatory.

R						$\pi_{\mathbf{a}, \mathbf{c}, \mathbf{e}}(R)$			S			$\sigma_{\mathbf{a} > 3 \mathbf{b} < 5}(S)$					
a	b	c	d	e	f	a	c	e	a	b	c	a	b	c			
1	2	3	4	5	6	1	3	5	3	2	3	3	2	3			
1	1	1	1	1	1	1	1	1	2	1	7				8	2	2
2	2	2	2	2	2	2	2	2	8	2	2						
1	2	3	4	5	8				1	2	9						

The syntax in use here for projection is $\pi_{\text{attributes}}(T)$, and likewise for selection it's $\sigma_{\text{condition}}(T)$. The former returns a relation with only the listed attributes, and the only takes rows which satisfy the given condition. Neither will return duplicate rows. Here is an example of a Cartesian product;

R		S			$R \times S$				
a	b	c	d	e	a	b	c	d	e
1	2	1	2	3	1	2	1	2	3
3	4	4	5	6	1	2	4	5	6
		7	8	9	1	2	7	8	9
					3	4	1	2	3
					3	4	4	5	6
					3	4	7	8	9

We can take a natural join, where the resulting relation contains all the tuples that have matching attributes in R , and S . This is less common, and we'd typically use something closer to $R \bowtie_{R.b=S.b} S$. Note that in this case, they are the same, since that's the only attribute that overlaps;

R		S			$R \bowtie S$			
a	b	b	c	d	a	b	c	d
1	2	1	2	3	1	2	5	6
1	3	2	5	6	3	4	8	9
3	4	4	8	9	3	4	9	9
		4	9	9				

Attributes can also be renamed with the notation $\rho_{\text{new/old},...}(R)$ notation, which is fairly self-explanatory.

Functional Dependencies

Relational database schemas should be normalised, which helps reduce redundancy, and avoids update, and deletion anomalies. The basis for normalisation is the concept of **functional dependency**, and **keys**. Consider the following example, which can experience update anomalies;

producer	city	product no.	price	quantity
3	London	52	65	4
22	Birmingham	10	15	5
22	Birmingham	12	4	11
3	London	44	43	32
3	London	43	3	27

If an existing producer changes address, and tuples aren't updated, then it leads to incoherence (as the data would now be invalid). If a new tuple / row is inserted, and the producer already exists, but with an older address, we will once again have incoherence. If a producer doesn't have any open orders, say rows 1, and 5 didn't exist, then data about producer 3 would be lost. The price of the product doesn't depend on the producer, and the location of the producer doesn't have anything to do with the product. There is also redundant data, since the address is duplicated. Ideally, the data should be stored like this;

Open Orders			Producer		Product	
producer no.	product no.	quantity	producer no.	city	product no.	price
3	52	4	3	London	52	65
22	10	5	22	Birmingham	10	15
22	12	11			12	4
3	44	32			44	43
3	43	27			43	3

A functional dependency is a constraint that if two tuples of a relation R agree on a set of attributes A_1, A_2, \dots, A_n , then they must also agree on the set of attributes B_1, B_2, \dots, B_m . This can be written as $A_1, A_2, \dots, A_n \rightarrow B_1, B_2, \dots, B_m$, and therefore the B set is functionally dependant on the A set, or the A set functionally determines the B set. Therefore, for any set of values A , then there is only one set of values B . The normal forms a designer can use are defined in terms of their functional dependencies. This can be summarised by the following examples;

functional dependency	for all tuple pairs, x, y , if	then assert
$A \rightarrow K$	$x.A = y.A$	$x.K = y.K$
$A, B \rightarrow K$	$x.A = y.A$ $x.B = y.B$	$x.K = y.K$
$A \rightarrow K, L$	$x.A = y.A$	$x.K = y.K$ $x.L = y.L$
$A, B \rightarrow K, L$	$x.A = y.A$ $x.B = y.B$	$x.K = y.K$ $x.L = y.L$

Note that if we have two functional dependencies, $A, B \rightarrow K$, and $A, B \rightarrow L$, it's trivial to combine them as $A, B \rightarrow K, L$ (or to split them). If we have the functional dependency $A, B, C, D, E \rightarrow A, C, X, Y, Z$, we can remove all the attributes on the right hand side, that are on the left hand side, such that the aforementioned FD reduces to $A, B, C, D, E \rightarrow X, Y, Z$. A **trivial** FD is one where all the attributes on the RHS are on the LHS.

Keys

If we have some set of attributes A_1, A_2, \dots, A_n , functionally determine all the remaining attributes of the relation, then the set A is a **superkey**. This also implies that two tuples cannot have the same superkey values, due to the uniqueness property. Let there be a set of all attributes in the relation, K , and the superkey set S . Therefore we can state $S \rightarrow \{x \in K \mid x \notin S\}$. A single relation may have more than one superkey, and superkeys can contain attributes that aren't strictly required. Generally, we are interested in superkeys where there isn't a subset of the superkey (hence the smallest superkey, by the irreducibility property), and this minimal superkey is referred to as the **candidate key**. If we have more than one candidate key, we can choose one of them to act as the primary key.

Closure

The set of all attributes functionally determined by a set L , under a set of dependencies F , the closure of L , denoted as L^+ . If L^+ contains all the attributes of R , then it follows that L is a superkey of R . If $\text{RHS} \subseteq \text{LHS}^+$, then $\text{LHS} \rightarrow \text{RHS}$ holds.

In order to compute the closure of a set of attributes, L , under a set of FDs given in the form $\text{LHS} \rightarrow \text{RHS}$, we keep adding all the LHSs which exist as a subset of L , and add the RHS to L . We repeat this until we can no longer add anymore RHSs. The final value is referred to as L^+ . Given the starting value of $L = \{A, B\}$, and the following FD set;

- (1) $A, B \rightarrow C$
- (2) $B, C \rightarrow A, D$
- (3) $D \rightarrow E$
- (4) $C, F \rightarrow B$

Starting with (1), both A , and B are in L , so we can add C . Now that C is in L , and B was already there, we can add D (no need to add A) from (2). Since we have D in L , we can add E . There is nothing else we can do, since F is not in L . Hence the value of $L^+ = \{A, B, C, D, E\}$.

Armstrong's Axioms

This is a sound (doesn't generate incorrect FDs), and complete (allows us to derive all valid FDs) axiomatisation of FDs. Given attributes A, B, C, \dots , and sets of attributes α, β, γ , we have the following axioms;

- reflexivity (trivial FDs) $\alpha \rightarrow \beta$ always holds if $\beta \subseteq \alpha$
- augmentation if $\alpha \rightarrow \beta$, then it follows $\alpha\gamma \rightarrow \beta\gamma$

- transitivity if $\alpha \rightarrow \beta$, and $\beta \rightarrow \gamma$, then $\alpha \rightarrow \gamma$
- (derived) union if $\alpha \rightarrow \beta$, and $\alpha \rightarrow \gamma$, then $\alpha \rightarrow \beta\gamma$
- (derived) decomposition if $\alpha \rightarrow \beta\gamma$, then $\alpha \rightarrow \beta$, and $\alpha \rightarrow \gamma$
- (derived) pseudotransitivity if $\alpha \rightarrow \beta$, and $\delta\beta \rightarrow \gamma$, then $\delta\alpha \rightarrow \gamma$

Closure, and Covers

Similar to attributes, we can find the closure of a FD set, which is the set of all FDs that can be inferred, denoted as F^+ . This can be approached by applying Armstrong's axioms. We start by initialising F^+ to the set F . By applying reflexivity, and augmentation, we can add new FDs to F^+ . Transitivity can then be applied to suitable FDs, and adding the derived ones to F^+ . This is repeated until no more changes occur.

We can say two FD sets F_1 , and F_2 are equivalent if they each imply the other - if any relation instance satisfying F_1 also satisfies F_2 , they are said to be covers of each other. A cover is canonical / minimal / irreducible if each LHS is unique, and no FD can be deleted from the cover, and still maintain an equivalent FD set. Additionally, no attributes can be deleted from any FD set, and still have an equivalent set. Therefore, the canonical cover has no redundant dependencies, nor attributes.

The method for testing whether some attribute, X , is extraneous depends on whether it's on the LHS, or RHS. For an attribute X on the LHS, we can say X is extraneous if $\text{RHS} \subseteq \{\text{LHS} - X\}^+$ under the FD set - note that here we are referring to closure on attributes, not FDs. On the other hand (no pun intended, I'm so depressed), an attribute X on the RHS is extraneous if $X \in \text{LHS}^+$ under the FD set, where X is removed from the RHS of all the FDs in the set.

To compute a canonical cover F for an FD set, we can use the following algorithm;

```

1 let  $F := \text{FD set}$ 
2 do
3     replace dependencies of  $\alpha \rightarrow \beta$  and  $\alpha \rightarrow \gamma$ , with  $\alpha \rightarrow \beta\gamma$ 
4     remove all extraneous attributes one at a time
5 until  $F$  doesn't change

```

Normalisation

As previously mentioned, in the previous section, normalisation reduces data duplication, and anomalies. One of the approaches taken to develop a normalised schema is to start with a few big relations, and then decompose them into a suitable normal form. The ones we'll look at are **Boyce-Codd Normal Form (BCNF)**, and **Third Normal Form (3NF)**. Consider the following example, which is just a "big" relation;

title	year	length	genre	studio	actor
Star Wars	1977	124	Science Fiction	Fox	Carrie Fisher
Star Wars	1977	124	Science Fiction	Fox	Mark Hamill
Star Wars	1977	124	Science Fiction	Fox	Harrison Ford
Gone with the Wind	1939	231	Drama	MGM	Vivien Leigh
Wayne's World	1992	95	Comedy	Paramount	Dana Carvey
Wayne's World	1992	95	Comedy	Paramount	Mile Myers

You'll notice straight away that there is a huge amount of redundant, repeated data, which is present in multiple rows. Updating the `length` of Star Wars would require 3 updates (for the relation to remain consistent), which is error prone. Inserting a new actor for a movie would require duplicating a fair bit of data, and would require checks to maintain consistency. Finally, deleting Vivien Leigh would delete the entire row, which deletes information about the movie.

Decomposition

Given a relation $R(A_1, A_2, \dots, A_n)$, we can decompose it into two projected relations S , and T , such that $\text{attr}(R) = \text{attr}(S) \cup \text{attr}(T)$. This also means that $S = \pi_{\text{attr}(S)}(R)$, and similar for T .

When we decompose a relation, it's crucial for us to ensure we can recover the original relation by joining the decomposed relations, and also preserve the FDs of the original relation. Note that the latter isn't always possible with BCNF.

Given R decomposed into S , and T , the decomposition is **lossless** if at least one of the following FDs hold in the closure of FD set of R ;

- $\text{attr}(S) \cap \text{attr}(T) \rightarrow \text{attr}(S)$
- $\text{attr}(S) \cap \text{attr}(T) \rightarrow \text{attr}(T)$
- this means the common attributes of S , and T form a superkey of either of the decomposed relations

If we're able to check the FDs of R without joining the decomposed sets, then the decomposition is **dependency preserving**.

Consider the following example, on $R(A, B, C)$, with the FD set $\{A \rightarrow B, B \rightarrow C\}$;

decomposition	lossless?	dependency preserving?
$S(A, B)$ $T(B, C)$	$B \rightarrow A, B$ doesn't hold $B \rightarrow B, C$ holds \therefore lossless	yes, $A \rightarrow B$ is checked with S , and $B \rightarrow C$ with T
$S(A, B)$ $T(A, C)$	$A \rightarrow A, B$ holds \therefore lossless $A \rightarrow A, C$ holds also	no, $B \rightarrow C$ cannot be checked without joining

Boyce-Codd Normal Form

A relation R is in BCNF, if and only if, for all non-trivial FDs (including derived FDs) of the relation, the LHS of every FD is a superkey (i.e. contains a key)

For example, if we have the FD $\text{title, year} \rightarrow \text{length, genre, studio}$, it's not in BCNF, as title, year doesn't functionally determine actor .