

# Fundamentos de Bases de Datos

Resumen práctico

Diego González  
2014  
diesgomo@gmail.com





# 1. Bases de Datos

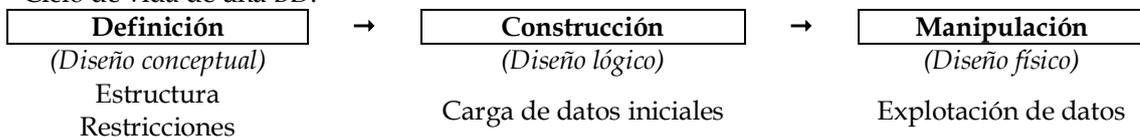
## 1.1 Definiciones Básicas

- **Dato:**  
Representación simbólica que ante la inexistencia de un contexto, carece de significado
- **Información:**  
Datos interpretados (que poseen una semántica)
- **Sistema de Información (SI):**  
Conjunto de componentes que interactúan con el objetivo de almacenar, recuperar y procesar datos e información con el fin de crear nueva información
- **Base de Datos (BD):**  
Conjunto de datos relacionados entre sí
- **Metadatos:**  
Reglas que determinan si una información pertenece a una BD
- **Sistema Manejador de Bases de Datos (DBMS):**  
SW que administra BDs
- **Modelo:**  
Lenguaje que permite describir:

Aspecto	Modelo	
Estructuras de Datos	De datos	Conceptual
Restricciones de integridad		
Operaciones (sobre la estructura)		

## 1.2 Bases de Datos

- Ciclo de vida de una BD:



- Ventajas de uso de BDs:

Aspecto	Descripción	Ventajas
Definición	De los datos abstractamente	<ul style="list-style-type: none"> <li>• Centralización</li> <li>• Abstracción (de programas-datos/operaciones)</li> <li>• Múltiples vistas</li> </ul>
Construcción	Programación de la BD	<ul style="list-style-type: none"> <li>• Eficiencia de consultas</li> <li>• Enfoque declarativo de restricciones e integridad</li> <li>• Estandarización<sup>1</sup></li> </ul>
Manipulación	Explotación de la BD	<ul style="list-style-type: none"> <li>• Datos compartidos</li> <li>• Control de concurrencia</li> <li>• Seguridad</li> <li>• Tolerancia ante fallas</li> <li>• Persistencia de datos consistentes (incluso de objetos)</li> </ul>

<sup>1</sup> De modelos y lenguajes que usan la BD

- Uso no adecuado de BDs:

Aspecto	Descripción
<b>Costo de inversión</b>	Alto en HW, SW y capacitación
<b>Costo de administración</b>	De la BD y el DBMS
<b>Naturaleza de los datos</b>	Pocos datos o muy estables en el tiempo
<b>Performance crítica</b>	Sistemas de tiempo real
<b>No hay concurrencia</b>	

- Componentes:

<b>Esquema</b>	Definición de la BD. Totalmente equivalente al Modelo Conceptual
<b>Instancia</b>	Estado de la BD y sus datos en un momento determinado

### 1.3 Actores

- Primarios:

Aquellos que interactúan tanto con la DB como con su información:

Actores	Funciones
<b>Administradores</b>	Autorizar, monitorear y coordinar
<b>Diseñadores</b>	Identificar los datos que se almacenarán y sus estructuras
<b>Desarrolladores</b>	Implementar los mecanismos de acceso
<b>Usuarios finales</b>	<ul style="list-style-type: none"> <li>• Casuales</li> <li>• Paramétricos: Utilizan transacciones enlatadas</li> <li>• Sofisticados</li> <li>• Independientes: Operan mediante interfaces sencillas</li> </ul>

- Secundarios:

Aquellos que no están interesados en el contenido de la BD:

Actores	Funciones
<b>Implementadores</b>	De módulos e interfaces de la DBMS
<b>Operadores de Sistemas</b>	Ejecutan las políticas definidas por los diseñadores

### 1.4 Modelos de datos

- Modelo de datos:

Colección de conceptos utilizados para describir la estructura de una BD

- Categorías:

Modelo	Descripción
<b>Conceptual</b>	Describe la relación y características de los datos
<b>Lógico</b>	Representa las características indicadas en el Modelo Conceptual. Ejemplos: Modelo Relacional                      Modelo Dato-Objeto
<b>Físico</b>	Describe cómo el modelo lógico es implementado en una estructura de archivos

- Arquitectura de tres esquemas:

#	Esquema	Descripción	Lenguaje	
			Acrónimo	Sigla
1	<b>Interno</b>	Almacena la estructura física de la BD	DDL	Data Definition Language
2	<b>Conceptual</b>	Describe la estructura de toda la BD y sus usuarios	SDL	Storage Definition Language
3	<b>Vistas de Usuario</b>	Describe determinada parte de la BD que es de interés para determinado grupo de usuarios	VDL	View Definition Language



- *Mapeo*:  
Acto de transformar y transmitir las consultas entre los esquemas de distintos niveles
- *Independencia de datos*:  
Habilidad de cambiar el esquema de un nivel sin afectar al esquema del nivel superior<sup>2</sup>
- *DML (Data Manipulation Language)*:  
Familia de lenguajes de manipulación de BD:

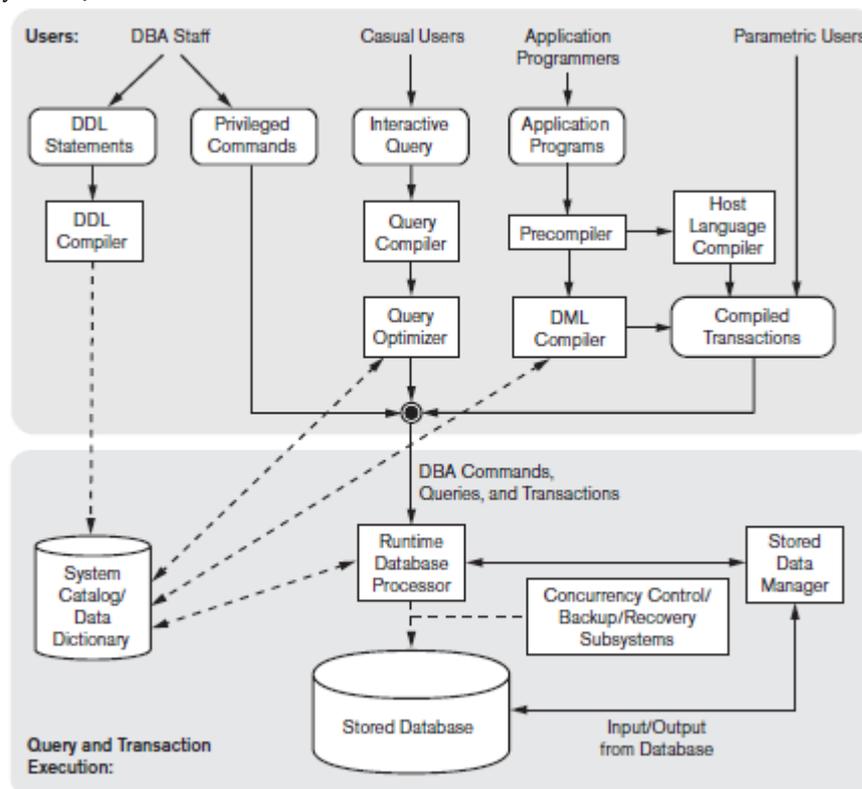
Categoría	Descripción
<b>Procedural</b>	Incrustado en lenguajes de programación de propósito general
<b>Consulta</b>	Utilizado de forma independiente e interactiva

- *SQL (Simply Query Language)*:  
Lenguaje que incluye a las familias DDL, SDL, VDL y DML

## 1.5 DBMS

- Interfaces de una DBMS:
  - Basada en Menús
  - Lenguaje Natural
  - Formulario
  - Entrada y Salida de Voz
  - Específicas del DBA
  - Gráficas
  - Usuarios Paramétricos

- Módulos y Componentes:



- Utilidades del DBMS:
  - Carga de archivos de datos
  - Copia de Seguridad de la BD
  - Organización y almacenamiento de la BD
  - Monitoreo de Performance

- Clasificación:

Criterio	Clasificación
<b>Modelo de Datos</b>	<ul style="list-style-type: none"> <li>• Relacional</li> <li>• Dato-Objeto</li> <li>• Jerárquica</li> <li>• Legacy</li> <li>• Objeto-Relacional</li> <li>• XML-nativo</li> </ul>
<b>Cantidad de Usuarios</b>	<ul style="list-style-type: none"> <li>• Único</li> <li>• Multiusuario</li> </ul>
<b>Distribución</b>	<ul style="list-style-type: none"> <li>• Centralizado</li> <li>• Distribuido<sup>3</sup></li> </ul>
<b>Tipo de Ruta de Acceso</b>	
<b>Especialización</b>	<ul style="list-style-type: none"> <li>• Propósito General</li> <li>• Propósito Específico</li> <li>• OLTP (Online Transaction Process)<sup>4</sup></li> </ul>

<sup>3</sup> Esta categoría, a su vez, puede clasificarse en **homogéneo** (utiliza el mismo SW en todas las locaciones) y **heterogéneo**

<sup>4</sup> Soporte a muchas transacciones concurrentes a bajo costo de demoras

## 2. Diseño Conceptual

### 2.1 Normas Básicas

- *Diseño conceptual:*  
Etapa en la que se construye un esquema conceptual en un lenguaje de alto nivel, definiendo el dominio del problema
- Construcción:
  1. Estudio de la realidad
  2. Especificación en lenguaje de alto nivel
  3. Validación del resultado

- Componentes:

Origen	Componente	Descripción
<b>Conceptos</b>	Conjuntos	Elementos de interés que por sus características, es conveniente la agrupación
	Relaciones	Entre conjuntos
	Restricciones de integridad	Validan los elementos que pueden pertenecer a una relación
<b>Términos</b>	Atributo	Características comunes de los elementos de un conjunto
	Cardinalidad	De elementos de un conjunto relacionados con el origen: N:1, N:N

- *Totalidad:*  
Una relación es total respecto de un conjunto, si todos sus elementos están en dicha relación

- Principios:

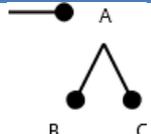
<b>100%</b>	Un esquema conceptual de un problema debe representar todos sus aspectos relevantes
<b>Conceptualización</b>	Un esquema conceptual no debe tener elementos de implementación

### 2.2 Modelo Entidad-Relación

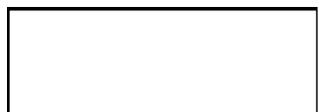
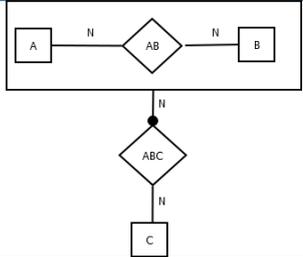
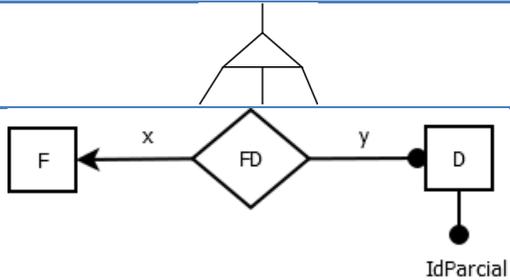
- *Modelo Entidad-Relación (ER):*  
DDL gráfico utilizado para la representación de estructura y restricciones de integridad<sup>5</sup>
- Construcción:
  1. Identificar elementos del problema
  2. Identificar relaciones entre los elementos
  3. Representar propiedades de los elementos
  4. Especificar restricciones

<sup>5</sup> Las mismas se realizan en lógica de primer orden o teoría de conjuntos

- Constructores:

Constructor	Interpretación
	<b>Conjunto de Entidades</b> <b>Atributos:</b> Función entidad-valor <sup>6</sup> : $A: \text{ConjuntoElementos} \rightarrow Bs \times Cs$ $B: A = Bs \times Cs \rightarrow Bs$
	<b>Atributo Multivalorado:</b> compuesto de muchos valores del mismo dominio
	<b>Atributo Determinante:</b> Actúa como "identificador" <sup>7</sup> , no permitiendo que dos entidades del conjunto posean el mismo valor en ese atributo <sup>8</sup>

- Operadores:

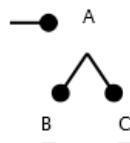
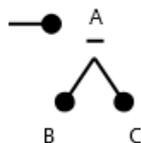
Operador	Interpretación				
	<b>Relación</b> entre dos o tres conjuntos de entidades (las relaciones triples pueden verse como la intersección de entidades en un Diagrama de Venn)				
	<b>Cardinalidad</b> (en el sentido de lectura): <table border="1"> <tr> <td><math>N</math></td> <td>Cualquier cantidad</td> </tr> <tr> <td><math>x</math></td> <td><math>x</math> cantidades como máximo</td> </tr> </table> <p>● Todo elemento debe estar en la relación</p>	$N$	Cualquier cantidad	$x$	$x$ cantidades como máximo
$N$	Cualquier cantidad				
$x$	$x$ cantidades como máximo				
	<b>Operador de "casting"</b> que reinterpreta su contenido como un nuevo conjunto de entidades. Ejemplo: 				
	<b>Especialización:</b> Indica que los "hijos" son subconjuntos <sup>9</sup> de los "padres" <b>Entidad Débil:</b> Para su identificación ("razón de ser") es necesario viajar hacia la entidad fuerte. Las entidades débiles pueden o no poseer un identificador parcial (que los identifique dentro de todas las relaciones con las entidades fuertes)				

<sup>6</sup> Para cada entidad, devuelve un valor. Pueden pensarse como "tipos abstractos"

<sup>7</sup> Estos son necesarios para cada conjunto de entidades: cada entidad debe ser identificable por un atributo determinante

<sup>8</sup> En particular, prestar atención en las combinaciones:

La combinación de  $B$  y  $C$  conforman un atributo "identificadorio"



No puede repetirse ningún  $B$  independientemente de que los  $C$  sean distintos y viceversa

<sup>9</sup> En el sentido matemático; por tanto, no se cumple si el cubrimiento del "padre" por sus "hijos", ni que los mismos sean mutu excluyentes

- Autorelaciones:  
Deben tener un rol en cada arco. En las restricciones se debe indicar si se desea (+) o no (-):
  - Ciclos
  - Reflexividad
  - + Transitividad

## 2.3 Calidad de los Esquemas Conceptuales

- Medidores de calidad:

Maximizar	Completitud
	Correctitud
	Minimalidad
Balancear	Expresividad
	Explicitud

- Completitud:  
Representa todas las características del problema
- Correctitud:
  - Sintáctica
  - Semántica
- Minimalidad:  
Cada elemento del problema aparece una sola vez en el esquema
- Expresividad:  
Facilidad de comprensión utilizando semántica del modelo
- Explicitud:  
No utiliza más formalismos que el diagrama ER



# 3. Modelo Relacional

## 3.1 Estructuras

- *Dominio:*  
Conjunto de valores atómicos
- *Relación:*

<b>Esquema</b>	$R(A_1, \dots, A_n)$	$R$ nombre de la relación $A_i$ atributo $i$ de dominio $D_i$
<b>Instancia</b>	$r(R)$	Conjunto <sup>10</sup> de tuplas que cumplen con el esquema de $R$
- *Tupla:*  
Función elemento de  $r(R)$  que, dado un atributo del esquema de  $R$ , retorna su valor
- *Esquema Relacional (Tablas):*  
Conjunto de esquemas de relación

## 3.2 Restricciones de Integridad (RI)

- Restricciones de Dominio:  
Tipado de cada Dominio
- Restricciones de Claves:
 

Tipo	Definición
<i>Superclave</i>	Subconjunto de atributos de una relación que es único para cada tupla posible
<i>Clave<sup>11</sup></i>	Superclave minimal
<i>Clave Foránea (FK)</i>	Atributo en una relación que es clave en otra relación (y cuyos valores se corresponden a ella)
- *RI:*  
Conjunto de Restricciones de Dominio y Restricciones de Claves
- Validez de una BD:  
Una BD es válida sii satisface todas sus RI
- Operaciones de modificación:

Operación	Sintaxis	Conflictos con RI
<b>Inserción</b>	INSERT <i>tupla</i> INTO $R$	
<b>Eliminación</b>	DELETE FROM $R$ WHERE <i>condición</i>	Cuando haya una FK sobre $R$
<b>Actualización</b>	UPDATE $R$ SET <i>atributo1</i> = <i>valor1</i> ... <i>atributoN</i> = <i>valorN</i> WHERE <i>condición</i>	Al intentar actualizar una clave o FK

<sup>10</sup> Por tanto no hay elementos ni ordenados ni repetidos

<sup>11</sup> Éstas se indican subrayando los atributos en el esquema

### 3.3 Cálculo Relacional de Tuplas: Sintaxis y Semántica

- Consultas:

Especificación de un conjunto de tuplas por comprensión sobre el universo de tuplas:

$$\{ \langle t_1, \dots, t_n \rangle / \varphi(x_1, \dots, x_n) \} \quad \left| \quad t_i = \begin{cases} x_i \cdot A_k & \text{con } \begin{cases} x_i \text{ variable libre de } \varphi \\ A_k \text{ atributo de la tupla de una tabla} \end{cases} \\ c_i & \text{constante} \end{cases} \right. \\ \varphi \text{ fórmula de lógica de primer orden tal que } FV(\varphi) = \{x_1, \dots, x_n\}$$

- Términos de construcción de  $\varphi$ :

$x_i$  variable                       $c_i$  constante (de algún Dominio)                       $x_i \cdot A_j$  con  $A_j$  atributo

- Fórmulas de construcción de  $\varphi$ :

Fórmula	Semántica
$t_i \text{ operador } t_j$ <sup>12</sup>	
$\text{tabla}_j(x_i)$	Verdadera sii $\exists x_i: x_i \in \text{tabla}_j$ en la BD
$(\varphi \text{ conector Lógico } \psi)$	Verdadera si lo es en el momento de la consulta
$(\neg \varphi)$	
$\exists x_i. \varphi$	Equivalente: $\exists t \in P. \varphi$ $\exists t. (P(t) \wedge \varphi)$
$\forall x_i. \varphi$	$\forall t \in P. \varphi$ $\forall t. (P(t) \rightarrow \varphi)$ $(\forall t \in P \wedge t.a = b). \varphi$ $\forall t. (P(t) \wedge t.a = b \rightarrow \varphi) \wedge \exists s. (P(s) \wedge s.a = b)$

### 3.4 Cálculo Relacional de Tuplas: Seguridad y Pragmática

- Fórmulas inseguras:

Fórmulas que podrían generar infinitos resultados

- Fórmulas seguras:

Una expresión es segura independientemente del dominio sii todos los alores de su resultado pertenecen al dominio de la expresión

- Criterios para fórmulas seguras:

$\bigvee \varphi_i$                       Todas las variables libres no negadas deben aparecer en cada  $\varphi_i$

$\bigwedge \varphi_i$                       Cada variable libre debe aparecer no negada en al menos una  $\varphi_i$

- Pragmática de las consultas:

Se accede a las claves de un elemento y luego se relacionan a través de ellas

### 3.5 Cálculo Relacional de Dominios y Equivalencias

- Consultas:  
Especificación de un conjunto de valores de dominio por comprensión sobre el universo de todos los dominios:

$$\{t_1, \dots, t_n / \varphi\} \quad \left| \quad t_i = \begin{cases} x_i \text{ variable libre de } \varphi \\ c_i \text{ constante} \end{cases} \right. \\ \varphi \text{ fórmula de lógica de primer orden tal que } FV(\varphi) = \{t_1, \dots, t_n\}$$

- Variables:  
Son atributos y se corresponden entre sí al especificar los esquemas y se utilizan en toda la estructura: se usa “\_” para omitir los valores en el mismo. Ejemplo:  
 $\{nom, dir / \exists nf. (FABS(nf, nom, dir) \wedge VENTAS(nf, \dots))\}$

- Equivalencias:

Tipo	Equivalencia	
Cuantificador	$\forall x. P(x)$	$\neg \exists x. \neg P(x)$
	$\exists x. P(x)$	$\neg \forall x. \neg P(x)$
Conector lógico	$\neg(x \wedge y)$	$\neg x \vee \neg y$
	$\neg(x \vee y)$	$\neg x \wedge \neg y$
	$x \rightarrow y$	$\neg x \vee y$

### 3.6 Álgebra Relacional

- Descripción de los operadores:  
Los operadores retornan siempre una relación. Cada relación -y por tanto composición de operadores- conforma una expresión. Como los resultados son conjuntos, no existen tuplas repetidas
- Definición de los operadores básicos:

Nombre	Notación	Descripción	Ejemplo
Relación	<i>Rel</i>	Expresión que retorna una copia de la relación <i>Rel</i>	<i>Student</i>
Selector	$\sigma_{cond} Expr$	Retorna las filas de la relación <i>Expr</i> que cumplen la condición lógica <i>cond</i>	$\sigma_{GPA > 3,5 \wedge sName = 'Fulano'} Student$
Proyector	$\pi_{A_{i_1}, \dots, A_{i_n}} Expr$	Retorna las columnas de atributos $A_{i_1}, \dots, A_{i_n}$ de la relación <i>Expr</i> <sup>13</sup>	$\pi_{SID, dec}(\sigma_{GPA > 3,5} Student)$
Producto Cartesiano (Cruz)	$Expr_1 \times Expr_2$	Retorna el producto cartesiano de ambas expresiones <sup>14</sup>	<i>Student</i> $\times$ <i>Apply</i>
Unión	$Expr_1 \cup Expr_2$	Retorna el conjunto unión de tuplas de sus operandos (el esquema debe ser el mismo o compatible <sup>15</sup> )	$\pi_{cName} College \cup \pi_{sName} Student$
Diferencia	$Expr_1 - Expr_2$	Retorna el conjunto diferencia de tuplas de sus operandos (el esquema debe ser el mismo o compatible)	$\pi_{SID} Student - \pi_{SID} Apply$

<sup>13</sup> Puede utilizarse posiciones de una relación mediante \$posición. Se deben utilizar todas posiciones o todos nombres de atributos

<sup>14</sup> Puede pensarse como la “unión” de las tablas que representan ambas Expresiones, donde los atributos de igual nombre son diferenciados:  $Expr_1.a \neq Expr_2.a$

<sup>15</sup> Por “compatible” se entiende distintos nombres de los atributos pero igual dominio



- Definición de operadores derivados:

Nombre	Notación	Definición	Descripción
Unión Natural	$Expr_1 \bowtie Expr_2 = Expr_1 * Expr_2$	$\left( \begin{array}{l} \pi_{Esquema(Expr_1) \cup Esquema(Expr_2)} \\ \left( \sigma_{\substack{Expr_1 A_1 = Expr_2 A_1 \wedge \\ Expr_1 A_2 = Expr_2 A_2 \wedge \dots}} (Expr_1 \times Expr_2) \right) \end{array} \right)$	Une las tuplas de sus operandos cuyos atributos de igual nombre adquieran igual valor <sup>16</sup>
Unión Theta	$Expr_1 \bowtie_{\theta} Expr_2$	$\sigma_{\theta}(Expr_1 \times Expr_2)$	Une sus operandos en la condición $\theta$
División	$Expr_1 \div Expr_2$	$\pi_X(Expr_1) - \pi_X((\pi_X(Expr_1) \times Expr_2) - Expr_1)$	Primeros $ X $ atributos de $Expr_1$ cuyos restantes atributos (coincidentes con $Expr_2$ ) incluyen todos los valores de $Expr_2$ <sup>17</sup>
Intersección	$Expr_1 \cap Expr_2$	$Expr_1 - (Expr_1 - Expr_2)$ Si el esquema coincide: $Expr_1 \bowtie Expr_2$	Intersecta sus operandos
Renombrador	$\rho_{Rel(A_{i_1}, \dots, A_{i_n})}(Expr)$	Retorna un nuevo esquema al computar $Expr$ creando una relación $Rel$ de atributos $A_{i_1}, \dots, A_{i_n}$ <sup>18</sup>	$\rho_{c(name)}(\pi_{cName} College)$ $\cup \rho_{c(name)}(\pi_{sName} Student)$

- Propiedades de los operadores:
  - $A \bowtie B = A \times B \Leftrightarrow esquema(A) \neq esquema(B)$
  - $A \times \emptyset = \emptyset$

<sup>16</sup> Por tanto, elimina las columnas con nombres repetidos

<sup>17</sup> Ejemplo:  $Q = R \div S$ :

$R(A, C)$		$S(C)$		$Q(C)$
$a_1 c_1$				
$a_1 c_2$				
$a_2 c_1$		$c_1$		
$a_2 c_2$	$\div$	$c_2$	$=$	$a_2$
$a_2 c_3$		$c_3$		
$a_2 c_4$				
$a_3 c_1$				
$a_3 c_3$				

<sup>18</sup> Su principal utilidad es la de unificar esquemas y así poder aplicar otros operadores



## 4. Equivalencia entre MER y Modelo Relacional

### 4.1 Dependencias de inclusión

- *Dependencias de Inclusión:*

Dependencias del modelo relacional que especifica que:

$$\prod_{a_1, \dots, a_n} A \subseteq \prod_{b_1, \dots, b_n} B$$

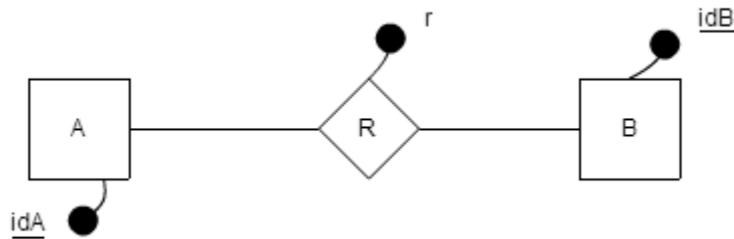
- Claves foráneas y dependencias de inclusión:  
Las claves foráneas son un caso particular de las dependencias de inclusión dado que en éstas no se exige la unicidad de aquella

### 4.2 Equivalencia de constructores

Constructor	Equivalencia
	ConjuntoDeEntidades(...)
	ConjuntoDeEntidades(..., B, C)
	ConjuntoDeEntidades( <i>idConjuntoDeEntidades</i> , ...) As( <i>idConjuntoDeEntidades</i> , A)
	Restricción: $\prod_{idConjuntoDeEntidades} As \subseteq \prod_{idConjuntoDeEntidades} ConjuntoDeEntidades$ ConjuntoDeEntidades( <u>A</u> , ...)

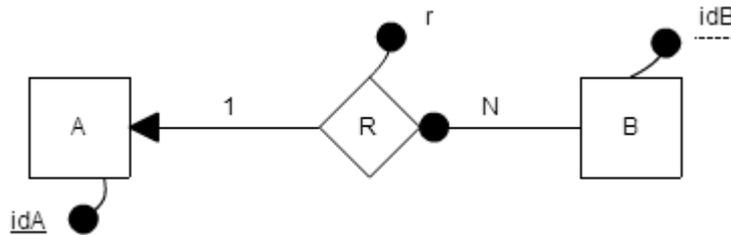
### 4.3 Equivalencia de relaciones<sup>19</sup>

- Binarias:



Cardinalidad (A : B)	Totalidad (de B a A)	Equivalencia	
N : N	×	$R(\underline{idA}, \underline{idB}, r, \dots)$ Restricciones: $\prod_{idA} R \subseteq \prod_{idA} A$	$A(\underline{idA}, \dots)$ $B(\underline{idB}, \dots)$ $\prod_{idB} R \subseteq \prod_{idB} B$
	✓	Se agrega la restricción:	$\prod_{idB} B \subseteq \prod_{idB} R$
1 : N	×	$R(\underline{idA}, \underline{idB}, r, \dots)$ Restricciones: $\prod_{idA} R \subseteq \prod_{idA} A$	$A(\underline{idA}, \dots)$ $B(\underline{idB}, \dots)$ $\prod_{idB} R \subseteq \prod_{idB} B$
	✓	Restricciones:	$A(\underline{idA}, \dots)$ $B(\underline{idB}, idA, r, \dots)$ $\prod_{idA} B \subseteq \prod_{idA} A$

- Entidades débiles:



Tablas:

$A(\underline{idA}, \dots)$

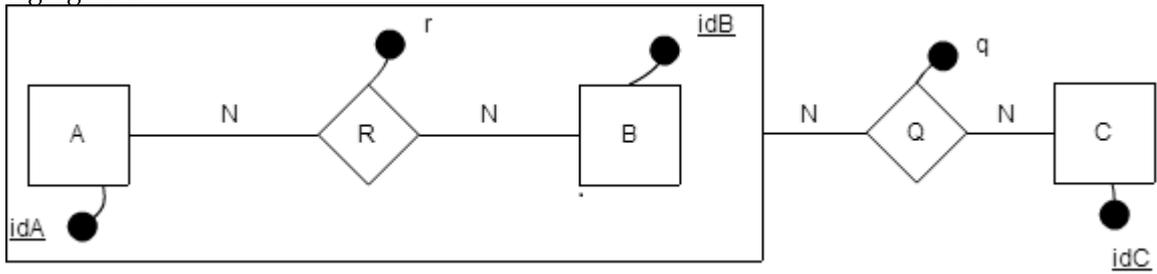
$B(\underline{idA}, \underline{idB}, r, \dots)$

Restricción:

$$\prod_{idA} B \subseteq \prod_{idA} A$$

<sup>19</sup> Los atributos son tratados como en 4.2

- Agregaciones<sup>20</sup>:



Tablas:

$C(\underline{idC}, \dots)$

$R(\underline{idA}, \underline{idB}, r, \dots)$

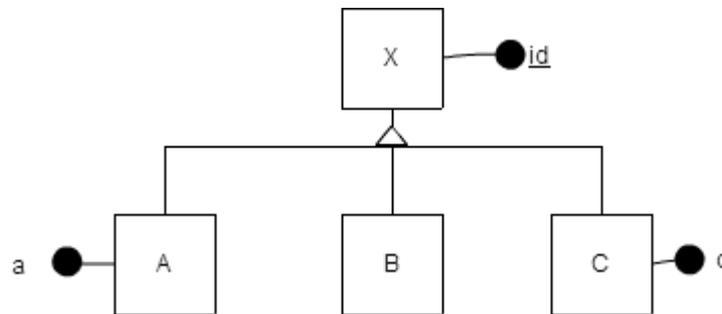
$Q(\underline{idA}, \underline{idB}, \underline{idC}, q, \dots)$

Restricciones:

$$\prod_{idC} Q \subseteq \prod_{idC} C$$

$$\prod_{idA, idB} Q \subseteq \prod_{idA, idB} R$$

#### 4.4 Equivalencia de categorizaciones



Restricciones	Equivalencia	
Ninguna <sup>21</sup>	$X(\underline{id}, \dots)$	$A(\underline{id}, a, \dots)$ $B(\underline{id})$ $C(\underline{id}, c, \dots)$
$X = A \cup B \cup C$	$X = \left\{ \langle t, id, \dots \rangle / \{A(t) \vee B(t) \vee C(t)\} \right\}$ <sup>22</sup>	$A(\underline{id}, a, \dots)$ $B(\underline{id})$ $C(\underline{id}, c, \dots)$
Disjuntas dos a dos	$X(\underline{id}, \dots, a, c, \text{tipo})$ Donde se incide en <i>tipo</i> a qué categoría pertenece	
Disjuntas (no dos a dos)	$X(\underline{id}, \dots, a, c, esA, esB, esC)$ Donde $esA, esB, esC$ son booleanos que identifican la agregación	

<sup>20</sup> Para otras cardinalidades se debe proceder como se indica anteriormente teniendo en cuenta que la agregación es una operación de casting

<sup>21</sup> Es aplicable a todos los demás casos con las restricciones adecuadas

<sup>22</sup> X se implementa como una vista



# 5. Diseño Relacional

## 5.1 Pautas de diseño

Esquema	Nombre	Descripción	Consideraciones
	Clara semántica de atributos	Claridad sobre qué representa una tabla y que una tupla esté en ella	No combinar atributos de varios tipos de entidades y vínculos en una sola relación
$\begin{matrix} a \\ b \\ \cancel{b} \end{matrix}$	Reducción de valores redundantes	Eliminar la mayor cantidad de valores redundantes	Inserción, eliminación y modificación anómalas
$\begin{matrix} a \\ \text{null} \\ b \end{matrix}$	Prevención de valores nulos	Evitar los valores nulos. Si no es posible, que su significado sea único	Que el significado de nulos sea: "no se conoce ese dato para esa tupla" <sup>23</sup>
$\begin{matrix} X(\underline{a}, \underline{b}, c, e) \\ Y(\underline{c}, \underline{e}) \\ Y(\underline{b}, c, e) \end{matrix}$	Particionamiento de tablas consistente	Tuplas erróneas: los esquemas deben reunirse por condición de igualdad sobre atributos claves	Sólo "dividir" las tablas sobre atributos claves

## 5.2 Dependencias

- Dependencia Funcional (DF) y Dependencias Multivaluadas (DMV): Restricciones de integridad sobre el modelo relacional que implica que:

Dada una instancia de una relación para la que X e Y son atributos y Z = R - XY:

Dependencia	Funcional	Multivaluada
<b>Sigla</b>	DF	DMV
<b>Notación</b>	$X \rightarrow Y$	$X \twoheadrightarrow Y$
<b>Definición</b>	$\forall t_1, t_2 \in r, t_1[X] = t_2[X] \Rightarrow t_1[Y] = t_2[Y]$	$\forall r, \text{ si } \exists t_1, t_2: t_1[X] = t_2[X] \Rightarrow \exists t_3, t_4 \text{ tales que:}$ <ul style="list-style-type: none"> <li><math>t_1[X] = t_2[X] = t_3[X] = t_4[X]</math></li> <li><math>t_3[Y] = t_1[Y] \text{ y } t_4[Y] = t_2[Y]</math></li> <li><math>t_3[Z] = t_2[Z] \text{ y } t_4[Z] = t_1[Z]</math></li> </ul>
<b>Descripción</b>	X determina a Y para toda tupla de esa instancia	X determina para cada Y distinto, un conjunto (el mismo y todo) de valores de Z <sup>24</sup>

- DMV Embebida: Una DMV Embebida  $X \twoheadrightarrow Y | Z$  indica que  $X \twoheadrightarrow Y$  se cumple para todo esquema  $X \cup Y \cup Z$ <sup>25</sup>

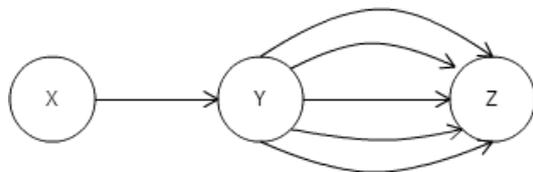
<sup>23</sup> Especial atención sobre categorizaciones

<sup>24</sup> Z es independiente de Y y viceversa, pero se relacionan en R por la existencia de X:

Por cada valor  $Z_i$  de  $X_0Y_0$ , si existe una tupla

$X_0Y_1$ , ésta debe estar relacionada también con  $Z_i$ :

Tupla	X	Y	Z
$t_1$	$x_1$	$y_1$	$z_1$
$t_2$	$x_1$	$y_2$	$z_2$
$t_3$	$x_1$	$y_1$	$z_2$
$t_4$	$x_1$	$y_2$	$z_1$



<sup>25</sup> Por tanto, durante una descomposición, éstas no deben considerarse hasta que se obtenga que algún  $R_i = X \cup Y \cup Z$



- Algoritmo de eliminación para encontrar  $F_{min}$ :

#	Algoritmo	Aceleración <sup>28</sup>	Descripción
1	$G = F$		Comenzar por dependencias existentes
2	Reemplazar $X \rightarrow AB \dots$ por $X \rightarrow A, X \rightarrow B, \dots$		"Partir" las DF
3	Para cada $X \rightarrow A \in G$ , para cada $B \in X$ , hallar $(X - B)_G^+$ Si $A \in (X - B)_G^+$ , reemplazar $X \rightarrow A$ por $X - B \rightarrow A$ en $G$	Atributos que sean compuestos a la izquierda, y se analiza si sacando uno de ellos se infiere el mismo resultado	Eliminar atributos redundantes
4	Para cada $X \rightarrow A \in G$ , hallar $(X)_G^+(X \rightarrow A)$ Si $A \in (X)_G^+(X \rightarrow A)$ , eliminar $X \rightarrow A$ en $G$	Analizar aquellas DF cuyo atributo a la derecha esté repetido en otra DF	Eliminar dependencias redundantes

### 5.3 Definiciones sobre Atributos y Dependencias

- Definiciones sobre claves:

Clave	Definición
<b>Superclave</b>	$S \subseteq R$ es una <i>superclave</i> sii $\forall r$ instancia de $R$ , $\nexists t_1, t_2 \in r$ tal que $t_1[S] = t_2[S]$
<b>Clave</b>	Una <i>clave</i> es una superclave minimal <sup>29</sup>
<b>Clave candidata</b>	Cada una de las claves existentes de una relación es una <i>clave candidata</i>
<b>Clave primaria</b>	Clave de una relación asignada arbitrariamente entre las claves candidatas presentes
<b>Clave secundaria</b>	Claves que no son primarias
<b>Atributo primo</b>	Atributo miembro de alguna clave

- Definiciones sobre dependencias:

Dependencia $X \rightarrow Y$	Definición
<b>Total</b>	DF para la que no es posible quitar $A \subseteq X$ sin que deje de ser una DF <sup>30</sup>
<b>Parcial</b>	DF no Total: $\exists A \subset X: A \rightarrow Y$
<b>Transitiva</b>	$\exists Z$ no subconjunto de clave y $X \rightarrow Z$ y $Z \rightarrow Y$ <sup>31</sup>

<sup>28</sup> Criterio para seleccionar rápidamente las DF que deben considerarse para estudiar en cada caso

<sup>29</sup> Esto es, si se le quita alguno de sus atributos, deja de ser una superclave

<sup>30</sup> Los atributos del antecedente ya son "minimales" en el sentido de la existencia de la DF

<sup>31</sup> Por tanto, puede expresarse de forma transitiva usando toras DFs.

- Definiciones sobre descomposiciones:

Descomposición D	Definición
D de R	$D = \{R_i\}_{i=1,\dots,m}$ tales que $\bigcup_{i=1}^m R_i = R$
Proyección de F sobre $R_i$ <sup>32</sup>	$\prod_{R_i} F = \{X \rightarrow Y \in F^+ : X \cup Y \subseteq R_i\}$
Condición de preservación de F en D	$\bigcup_{i=1}^m \prod_{R_i} F = F^+$
D que cumple Join sin pérdida (JSP)	<p>Si <math>i = 2</math><sup>33</sup> debe cumplirse una de las siguientes condiciones<sup>34</sup>:</p> <p>i. <math>(R_1 \cap R_2) \rightarrow (R_1 - R_2) \in F^+</math></p> <p>ii. <math>(R_1 \cap R_2) \rightarrow (R_2 - R_1) \in F^+</math></p> <p>Si <math>i &gt; 2</math>: Para cada instancia <math>r</math> de <math>R</math> se cumple que</p> $\bigotimes_{i=1}^m \prod_{R_i} (r) = r$

- Cálculo de claves:

- Se definen tres conjuntos:

#	Conjunto de Atributos	Propiedad
1	No determinados por nadie	Nunca están a la derecha
2	No forman parte de ninguna clave	Sólo están a la derecha
3	Demás atributos	Están tanto a la derecha como a la izquierda

- Se determinan todas las combinaciones de atributos de los conjuntos 1 y 3
- Se estudia la clausura de cada elemento de 2 para ver si incluyen a todos los atributos

- Test de JSP para  $i > 2$ :

- Crear  $M \in \mathcal{M}_{D \times \{A \text{ atributo de } R\}}$
- Para cada fila  $i$ , si:

$$A_j \in R_i \Rightarrow m_{ij} = a_j \text{ valor fijo}$$

$$A_j \notin R_i \Rightarrow m_{ij} = b_{ij} \text{ variable del algoritmo}$$

- Repetir hasta que:
  - No haya modificaciones en  $M$  ó
  - Exista una fila con todos los símbolos  $a$  (sin importar el subíndice)

Lo siguiente:

Para cada  $X \rightarrow Y \in F$  se "fuerza" la dependencia:

En las filas donde los atributos de  $X$  coincidan

Igualar los símbolos de los atributos de  $Y$

(Como los  $b_{ij}$  son variables, si un  $b_{i_0 j_0}$  adquiere un valor  $a_k$ , todas sus referencias deberán hacerlo)

- La descomposición es JSP sii existe al menos una fila con todos los símbolos  $a$  (sin importar el subíndice)

<sup>32</sup> Dependencias de  $F^+$  donde todos los atributos que aparecen están en  $R_i$

<sup>33</sup> La definición si  $i > 2$  incluye si  $i = 2$ , pero se discrimina en este documento porque la definición para  $i = 2$  presenta muchas conveniencias prácticas

<sup>34</sup> Esto es, la intersección de los esquemas de la descomposición determina la diferencia de los esquemas en algún orden. Para las DMV, la definición es análoga reemplazando  $\rightarrow$  por  $\rightarrow$

Ejemplo:

$R = \{NSS, NOMBREE, NÚMEROP, NOMBREPR, LUGARP, HORAS\}$        $D = \{R1, R2, R3\}$   
 $R1 = EMP = \{NSS, NOMBREE\}$   
 $R2 = PROYECTO = \{NÚMEROP, NOMBREPR, LUGARP\}$   
 $R3 = TRABAJA_EN = \{NSS, NÚMEROP, HORAS\}$

$F = \{NSS \rightarrow NOMBREE; NÚMEROP \rightarrow \{NOMBREPR, LUGARP\}; \{NSS, NÚMEROP\} \rightarrow HORAS\}$

	NSS	NOMBREE	NÚMEROP	NOMBREPR	LUGARP	HORAS
R1	a <sub>1</sub>	a <sub>2</sub>	b <sub>13</sub>	b <sub>14</sub>	b <sub>15</sub>	b <sub>16</sub>
R2	b <sub>21</sub>	b <sub>22</sub>	a <sub>3</sub>	a <sub>4</sub>	a <sub>5</sub>	b <sub>26</sub>
R3	a <sub>1</sub>	b <sub>32</sub>	a <sub>3</sub>	b <sub>34</sub>	b <sub>35</sub>	a <sub>6</sub>

(matriz original S al principio del algoritmo)

	NSS	NOMBREE	NÚMEROP	NOMBREPR	LUGARP	HORAS
R1	a <sub>1</sub>	a <sub>2</sub>	b <sub>13</sub>	b <sub>14</sub>	b <sub>15</sub>	b <sub>16</sub>
R2	b <sub>21</sub>	b <sub>22</sub>	a <sub>3</sub>	a <sub>4</sub>	a <sub>5</sub>	b <sub>26</sub>
R3	a <sub>1</sub>	<del>b<sub>32</sub></del> a <sub>2</sub>	a <sub>3</sub>	<del>b<sub>34</sub></del> a <sub>4</sub>	<del>b<sub>35</sub></del> a <sub>5</sub>	a <sub>6</sub>

(la matriz S después de aplicar las dos primeras dependencias funcionales - la última fila sólo tiene símbolos "a", así que nos detenemos)

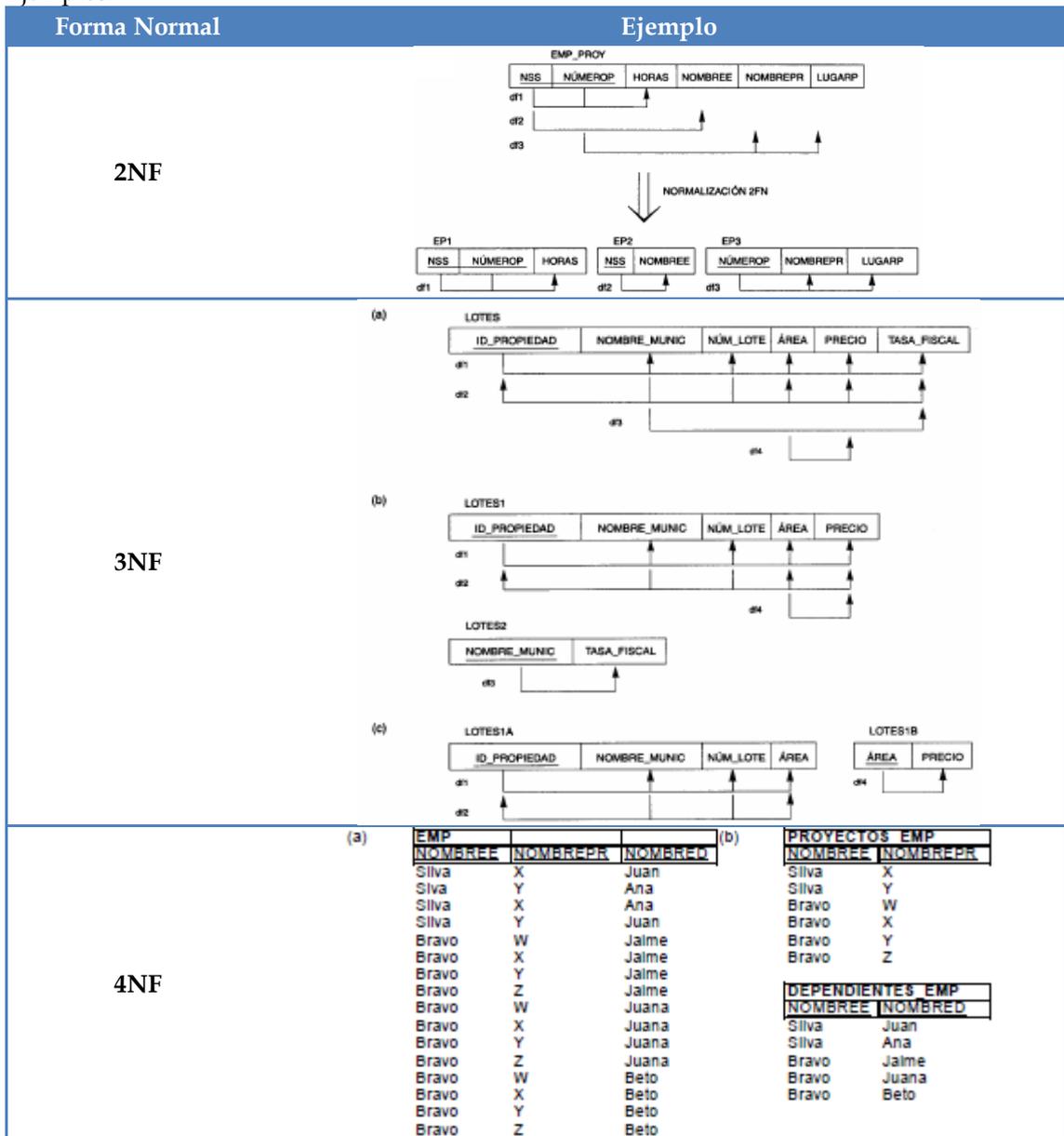
## 5.4 Formas Normales

- *Forma Normal:*  
Proceso de normalización permite:
  - Catalogar los esquemas
  - Construir un esquema normalizado a partir de uno no normalizado
- Formas Normales:

Abreviación	Forma Normal	Definición	Descripción
1NF	Primera	Los dominios de los atributos incluyen sólo valores atómicos	Todo diseño relacional visto en el curso lo cumple
2NF	Segunda	Ningún atributo no primo depende parcialmente de una clave	Ningún atributo que está sólo a la derecha depende de una parte de una clave. Por tanto, se podría normalizar creando varias tablas donde esas partes de las que dependen sean las claves de los atributos que dependen de él únicamente
3NF	Tercera	$\forall X \rightarrow A \in R$ se tiene que se cumple una de las siguientes propiedades: <ul style="list-style-type: none"> <li>• X superclave</li> <li>• A atributo primo</li> </ul>	Está en 2NF y ningún atributo no primo depende transitivamente de una clave
BCNF	Boyce-Codd	$\forall X \rightarrow A \in R$ se tiene que X es una superclave	Todas las DF son una superclave
4NF	Cuarta	$\forall X \rightarrow Y$ y $X \twoheadrightarrow Y$ no trivial, X es una superclave	Está en BCNF y todas las DMV no triviales son una superclave



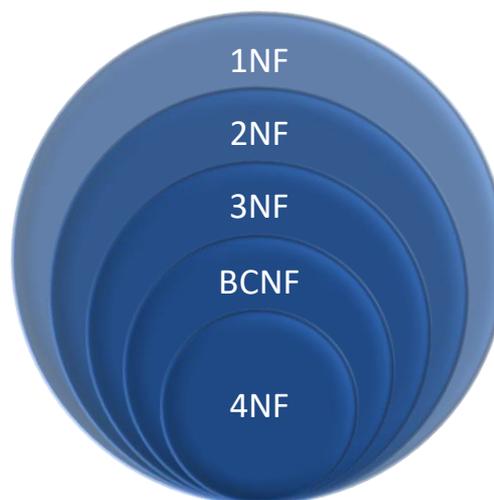
- Ejemplos:



- Algoritmos de descomposición de Formas Normales con JSP:

Forma Normal	Conserva DFs	Algoritmo
3NF	✓	<ol style="list-style-type: none"> <li>Hallar <math>F_{min}</math></li> <li>Hacer una tabla para cada <math>f \in F_{min}</math> donde el lado izquierdo sea el mismo: Para cada <math>X</math> lado izquierdo de <math>f</math> Crear <math>\{X \cup A_1 \cup \dots \cup A_m\} \in D</math> Donde <math>X \rightarrow A_i \forall i = 1, \dots, m</math> sean las únicas <math>f</math> con <math>X</math> como miembro izquierdo</li> <li>Colocar los atributos restantes y <math>X</math> en una sola tabla</li> <li>Si ninguna tabla contiene una clave de <math>R^{35}</math>: Crear una tabla adicional que contengan los atributos que forman una clave<sup>36</sup> de <math>R</math></li> </ol>
BCNF <sup>37</sup>	×	<ol style="list-style-type: none"> <li>Hacer <math>D = \{R\}</math></li> <li>Para cada <math>d \in D</math> que no esté en BCNF: <ul style="list-style-type: none"> <li>Encontrar <math>X \rightarrow Y \in d</math> que, proyectándolo, se determine que viola BCNF</li> <li>Reemplazar <math>d</math> por los dos esquemas siguientes: <ul style="list-style-type: none"> <li><math>d - Y</math></li> <li><math>X \cup Y</math></li> </ul> </li> </ul> </li> </ol>
4NF	×	Análogo a BCNF reemplazando donde diga "BCNF" por "4NF"

- 4NF y MER:  
En general, al pasar de un MER al modelo relacional, se obtendrá que éste estará en 4NF; cuando esto no sucede, existen altas probabilidades de que el MER esté mal construido
- Inclusiones de las Formas Normales:



<sup>35</sup> Para esto, se calcula si desde cada partición, la clausura cubre a  $F^+$

<sup>36</sup> Para ello se aplica el algoritmo de cálculo de claves

<sup>37</sup> Para su ejecución en papel, es útil diagramarlo en forma de árbol e ir descomponiendo los  $d$  que violen BCNF. La unión de las hojas será el resultado buscado



# 6. Procesamiento y Optimización de Consultas

## 6.1 Índices

- Organización física de los datos por un DBMS:  
Almacenan sus datos en un gran archivo o una partición dedicada, administrándolos en primera persona y utilizando los servicios más básicos del SO o incluso implementando los suyos propios<sup>38</sup>.

- *Índice*:  
Estructura de índice que, dada su clave, retorna como valor un puntero a un conjunto de bloques<sup>39</sup> que contiene los datos indizados

- Tipos de índices:

Relación	Clasificación	Dominio
<b>Abstracción</b>	Físico	Bloque de disco
	Lógico	Otros índices
<b>Orden<sup>40</sup></b>	Ordenados	Datos ordenados
	No Ordenados	Datos no ordenados
<b>Densidad</b>	Densos	Todas las claves
	No Densos	Algunos valores
<b>Niveles</b>	Simple	
	Multiniveles	
<b>Temporalidad</b>	Permanentes	
	Auxiliares	

- Desventajas en el uso de índices:  
Elevado costo de mantención (para cada inserción, eliminación o actualización de las tuplas).

<sup>38</sup> Con el objetivo de acceder lo más rápido posible a distintas secciones de los archivos almacenados

<sup>39</sup> Este conjunto de todos los bloques conforman el archivo que representa la tabla, un índice apunta a una porción de esos bloques donde justo se encuentran los datos de interés

<sup>40</sup> No se refiere al orden de la estructura de índices, sino de los datos indizados

- Orden de índices:

Clasificación	Índice	Clave	Valor	Ejemplo																																				
Ordenados	Primario	Primaria	Bloque + Desplazamiento	<table border="1"> <thead> <tr> <th>nombre</th> <th>edad</th> <th>salario</th> <th>depto</th> </tr> </thead> <tbody> <tr><td>Alberto</td><td>28</td><td>10000</td><td>Depto1</td></tr> <tr><td>Ana</td><td>30</td><td>12000</td><td>Depto2</td></tr> <tr><td>Juan</td><td>35</td><td>12000</td><td>Depto3</td></tr> <tr><td>...</td><td>...</td><td>...</td><td>...</td></tr> <tr><td>Lucía</td><td>40</td><td>20000</td><td>Depto3</td></tr> <tr><td>Luis</td><td>34</td><td>12000</td><td>Depto1</td></tr> <tr><td>...</td><td>...</td><td>...</td><td>...</td></tr> <tr><td>Sandra</td><td>30</td><td>15000</td><td>Depto2</td></tr> </tbody> </table>	nombre	edad	salario	depto	Alberto	28	10000	Depto1	Ana	30	12000	Depto2	Juan	35	12000	Depto3	...	...	...	...	Lucía	40	20000	Depto3	Luis	34	12000	Depto1	...	...	...	...	Sandra	30	15000	Depto2
	nombre	edad	salario	depto																																				
Alberto	28	10000	Depto1																																					
Ana	30	12000	Depto2																																					
Juan	35	12000	Depto3																																					
...	...	...	...																																					
Lucía	40	20000	Depto3																																					
Luis	34	12000	Depto1																																					
...	...	...	...																																					
Sandra	30	15000	Depto2																																					
Clúster <sup>41</sup>	Atributo que no es clave primaria	Bloque + Desplazamiento del primer registro <sup>42</sup>	<table border="1"> <thead> <tr> <th>Nombre</th> <th>Edad</th> <th>salario</th> <th>depto</th> </tr> </thead> <tbody> <tr><td>Alberto</td><td>28</td><td>10000</td><td>Depto1</td></tr> <tr><td>Ana</td><td>30</td><td>12000</td><td>Depto2</td></tr> <tr><td>Sandra</td><td>30</td><td>15000</td><td>Depto2</td></tr> <tr><td>Luis</td><td>34</td><td>12000</td><td>Depto1</td></tr> <tr><td>Juan</td><td>35</td><td>12000</td><td>Depto3</td></tr> <tr><td>Lucía</td><td>40</td><td>20000</td><td>Depto3</td></tr> <tr><td>...</td><td>...</td><td>...</td><td>...</td></tr> </tbody> </table>	Nombre	Edad	salario	depto	Alberto	28	10000	Depto1	Ana	30	12000	Depto2	Sandra	30	15000	Depto2	Luis	34	12000	Depto1	Juan	35	12000	Depto3	Lucía	40	20000	Depto3	...	...	...	...					
Nombre	Edad	salario	depto																																					
Alberto	28	10000	Depto1																																					
Ana	30	12000	Depto2																																					
Sandra	30	15000	Depto2																																					
Luis	34	12000	Depto1																																					
Juan	35	12000	Depto3																																					
Lucía	40	20000	Depto3																																					
...	...	...	...																																					
No Ordenados	Secundario	Atributos ordenados	Lista de bloques que contienen el valor																																					

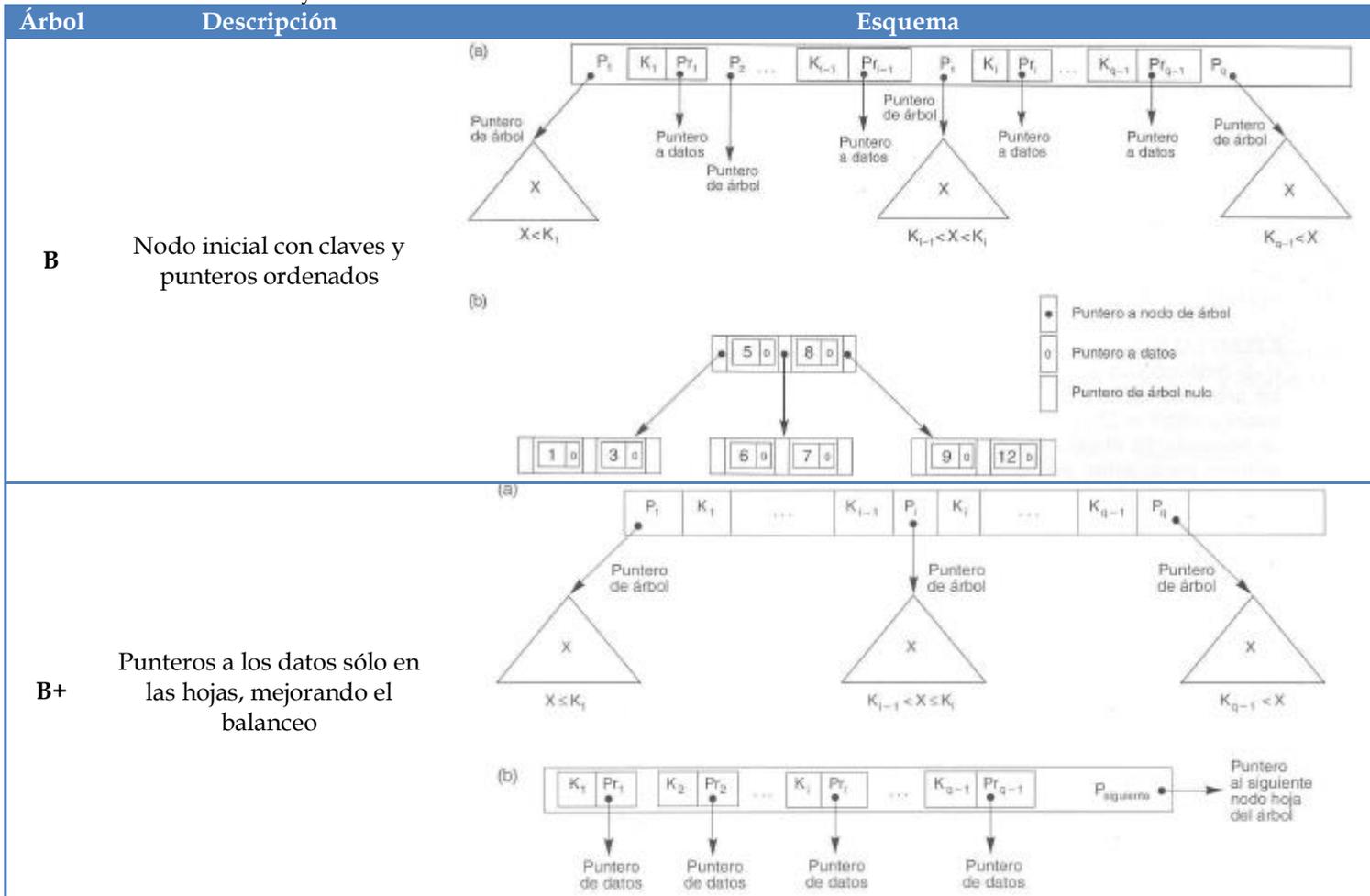
- *Bitmap*: Arreglo de bits sobre los valores de un atributo o una condición: para cada posición almacena si la tupla tiene determinado valor o cumple una condición predefinida para esa estructura
- Estructuras implementadoras de índices:

Estructura	Ventajas	Desventajas
<b>Hash</b>	Inserción y recuperación por igualdad	Relaciones de Orden
<b>Árboles</b>	Inserción y recuperación por igualdad y orden	Ocupa más espacio en disco
<b>Bitmap</b>	Atributos numerados	Pocas veces aplicable

<sup>41</sup> De agrupamiento

<sup>42</sup> Por tanto, para este caso debe verificarse que se encontraron todos los valores recorriendo a partir del obtenido por el índice; por este motivo, en general el valor es una lista de punteros

- Árboles B y B+:



- Comparación entre estructuras:  
Son comparadas utilizando la cantidad de acceso a disco teniendo en cuenta que en un acceso a disco podrían recuperarse más de un nodo de la estructura y que estos podrían buffearse en memoria

## 6.2 Proceso de Optimización

- Tipos de Optimización:

Tipo	Plan	Descripción
<b>Heurística</b>	Lógico	Reduce el árbol de consultas en álgebra relacional
<b>Costos</b>	Físico	A cada operación de los Planes Lógicos le asocia una o más implementaciones según las estructuras disponibles
	Final	Evalúa los Planes Físicos según los accesos de E/S requeridos

- Proceso de Optimización:

#	Etapas	Descripción	Ejemplo															
1	<b>Álgebra Relacional</b>	<p>Se construye la consulta en Álgebra Relacional partiendo de SQL como:</p> $\prod_{select} \sigma_{where}(x_{from})$ <p>Obteniendo un <i>árbol canónico</i></p>	<pre>select e.nombre, d.piso from departamentos d, empleados e where e.depto = d.nroD and e.salario &gt; 30000</pre>															
2	<b>Planes Lógicos</b>	<p>A partir del árbol canónico se reduce el árbol empleando:</p> <ul style="list-style-type: none"> <li>- <i>Equivalencias de expresiones</i></li> <li>- Catálogo del a DB</li> </ul> <p>En cada arista se indica la cantidad de tuplas estimadas que cumplen con las condiciones de la operación de destino</p>	<p>Catálogo de la DB:</p> <table border="1"> <thead> <tr> <th>Parámetro</th> <th>Valor</th> </tr> </thead> <tbody> <tr> <td>Tamaño de EMPLEADOS (tuplas)</td> <td>100.000</td> </tr> <tr> <td>Tamaño de DEPARTAMENTOS (tuplas)</td> <td>100</td> </tr> <tr> <td>Selectividad de <math>\sigma_{(salario &gt; 30000)}</math></td> <td>1/10.000</td> </tr> </tbody> </table> <p>Planes Lógicos:</p>	Parámetro	Valor	Tamaño de EMPLEADOS (tuplas)	100.000	Tamaño de DEPARTAMENTOS (tuplas)	100	Selectividad de $\sigma_{(salario > 30000)}$	1/10.000							
Parámetro	Valor																	
Tamaño de EMPLEADOS (tuplas)	100.000																	
Tamaño de DEPARTAMENTOS (tuplas)	100																	
Selectividad de $\sigma_{(salario > 30000)}$	1/10.000																	
3	<b>Planes Físicos</b>	Se consideran las implementaciones de $\sigma$ y $\bowtie$ para determinar los costos de cada plan	<table border="1"> <thead> <tr> <th>Operación</th> <th colspan="3">Implementación</th> </tr> </thead> <tbody> <tr> <td rowspan="2"><math>\sigma</math></td> <td>Búsqueda Lineal</td> <td>Búsqueda Binaria</td> <td rowspan="2">Índice</td> </tr> <tr> <td><math>\bowtie</math></td> <td>Loop anidado</td> <td>Loop único</td> </tr> <tr> <td></td> <td></td> <td></td> <td>Sort Merge</td> </tr> </tbody> </table>	Operación	Implementación			$\sigma$	Búsqueda Lineal	Búsqueda Binaria	Índice	$\bowtie$	Loop anidado	Loop único				Sort Merge
Operación	Implementación																	
$\sigma$	Búsqueda Lineal	Búsqueda Binaria	Índice															
	$\bowtie$	Loop anidado		Loop único														
			Sort Merge															
4	<b>Plan Final</b>	Calculados los costos de cada plan, se selecciona el de menor valor																

- Índices en el presente modelo teórico:

En el modelo teórico que aquí se construye, asume que los índices sólo pueden ser empleados en el primer nivel de las operaciones (según el Plan Lógico)

### 6.3 Optimización por Heurísticas

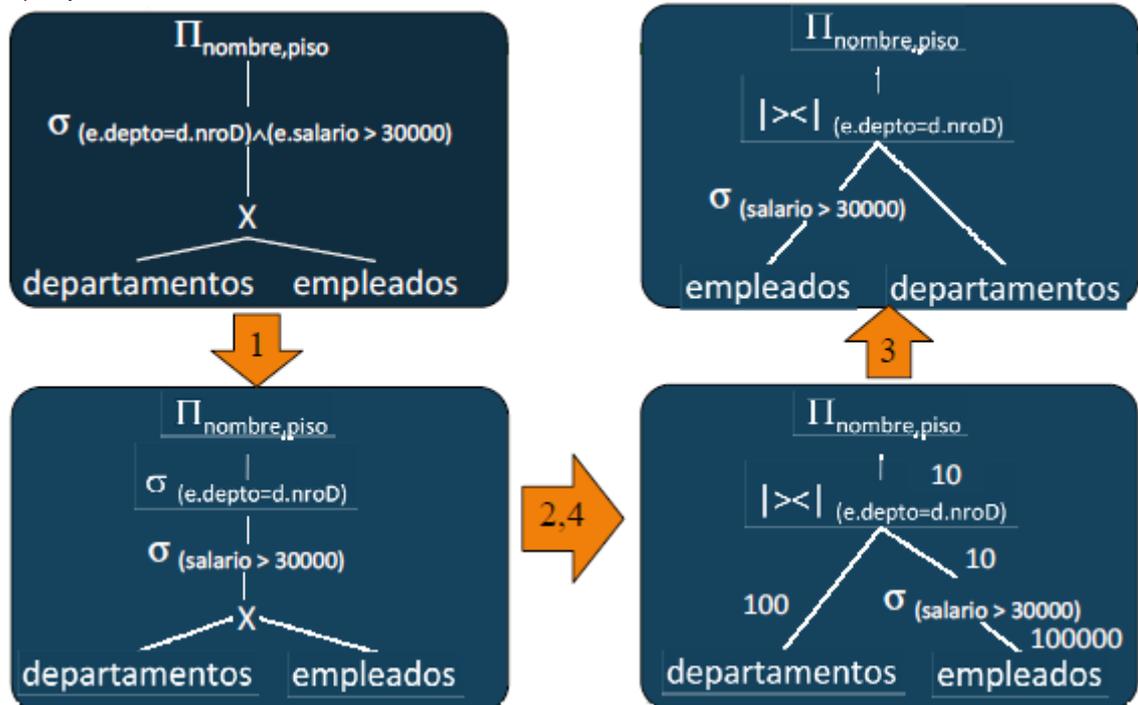
- Criterios:

#	Elemento	Acción
0	En todos los casos y siempre que sea posible, deben aplicarse las Equivalencias de Expresiones a fin de minimizar las operaciones de la consulta	
1	$\sigma$	Descomponer en $\sigma$ s más simples
2	$\sigma$	Mover lo más abajo posible
3		Mover a la izquierda las ramas con menos tuplas
4	$(\sigma, \times)$	Reemplazar por $\bowtie_{condición}$
5	$\pi$	Mover lo más abajo posible agregando los $\pi$ s necesarios <sup>43</sup>

- Equivalencia de expresiones:

Categoría	Equivalencia	Descripción/Aclaraciones
$\sigma$	Descomposición	
	Conmutatividad	
$\pi$	$\pi_{a_i}(\pi_{a_1, \dots, a_n} R) = \pi_{a_i} R$	Eliminación de redundancia
	$\pi_a \sigma_c R = \sigma_c \pi_a R$	Conmutación con $\sigma$ si éste contiene atributos de $\pi$
$\bowtie$	$\sigma_c(R \times Q) = R \bowtie_c Q$	Equivalencia de $(\sigma, \times)$
	Asociatividad de $\bowtie_c$	
	Conmutatividad de $\bowtie_c$	
$\times$	$\sigma_c(R \times Q) = (\sigma_c R) \times Q$	Si $c$ sólo contiene atributos de $R$
	$\pi_{a \cup b}(R \times Q) = \pi_a R \times \pi_b Q$	Si $a$ es de $R$ y $b$ es de $Q$
$\cup/\cap$	Reglas de conjuntos	
	$\sigma_c(R \cup Q) = \sigma_c R \cup \sigma_c Q$	Análogo para $\cap$ , $-$ y $\pi$ en lugar de $\sigma$

- Ejemplo:



<sup>43</sup> En general, no se aplica debido a la pérdida de índices y la necesidad de mantener varios atributos de proyección en los niveles inferiores porque serán utilizados en los niveles superiores

## 6.4 Optimización por Costos

- Parámetros:

Notación	Unidad	Definición	Fórmula
$n_T$	tuplas	...de la tabla $T$	
$R_T$	bytes	...de una tupla de $T$	
$b_T$	bloques	...necesarios para almacenar las tuplas de $T$	$\left\lceil \frac{n_T}{bfr_T} \right\rceil$
$bf_T$	tuplas	...en un bloque de $T$	$\left\lfloor \frac{\text{bytes del bloque}}{\text{bytes de la tupla}} \right\rfloor$
$x_T$	niveles	...de un índice de $T$	$\log_k n_T + 1^{44}$
$V(A, T)$		...valores distintos que tiene un atributo $A$ sobre una tabla $T$	$n_T^{45}$
$sl(\sigma_c T)$	tuplas	Fracción de tuplas que cumplen $c$ en $T$ original	$\frac{1}{V(A, T)^{46}}$
$js(R \bowtie_c S)$	tuplas	...que deben seleccionarse respecto del producto de $\times$	$\frac{1}{\min\{V(A, R), V(A, S)\}^{47}}$
$T(\sigma_c R)$	tuplas	...que cumplen con $c$ en $T$	$n_R \times sl(\sigma_c R)$
$T(R \bowtie_c S)$	tuplas	...que cumplen con $\bowtie_c$	$n_R \times n_S \times js(R \bowtie_c S)$

- Estrategia de implementación de operadores:

Pipelined	Ejecución paralela de operadores	Generación de resultados	Aplicación
✓	✓	Son comunicados entre operadores sin grabarse a disco <sup>48</sup>	$\pi$
✗	✗	Son comunicados entre operadores grabando resultados intermedios	$\sigma$ $\bowtie$

- Estimaciones de costos:

Los costos son calculados para cada bloque de disco<sup>49</sup> en las acciones de:

<b>Lectura</b>	Depende de la organización de los datos
<b>Escritura</b>	Costo de grabar <i>todo</i> el resultado $R$ : $\left\lceil \frac{n_R}{bfr} \right\rceil$

<sup>44</sup> Para un B+ con  $k$  punteros por nodo sobre clave

<sup>45</sup> Para un atributo clave

<sup>46</sup> Si la condición es una igualdad sobre  $A$  y se asume distribución uniforme

<sup>47</sup> Si es  $R \bowtie S$  natural sobre  $A$

<sup>48</sup> No hay costo de lectura ya que los datos están en el buffer

<sup>49</sup> Que por tanto pueden contener varios registros o índices

- Implementación de operadores:

Operador	Estrategia	Restricción	Descripción	Costo de lectura <sup>50</sup>									
$\sigma_c R$	Búsqueda Lineal		...de los registros que cumplen $c$	Caso promedio: $\frac{b_R}{2}$ Peor caso: $b_R$									
	Búsqueda Binaria	Ordenación	Se divide la cantidad de registro y se lee en el sentido de la ordenación	$\log_2 b_R + \left\lceil \frac{s}{b_R} \right\rceil - 1$									
	Índice Primario	Ordenación	Se consulta el índice	$x + 1$									
	Índice de Cluster	Ordenación	Se consulta el índice	$x + \left\lceil \frac{s}{bf_R} \right\rceil$									
	Índice Hash	$c$ relación de igualdad	Lectura sobre el índice	1 o 2									
	Índice secundario con B+		Se sigue la estructura B+	Peor caso: $x + s$									
$R \bowtie_c S$	Loop anidado por registros		Leída una tupla de $R$ , se leen todas las tuplas de $S$ 	$b_R + n_R \times b_S$									
	Loop anidado por bloques		Leído un bloque de $R$ , se leen todos los bloques de $S$ 	$b_R + \left\lceil \frac{b_R}{M-2} \right\rceil \times b_S^{51}$									
	Sort-Merge Join	Ordenación en $R$ y $S$	Recorrer $S$ y $R$ en paralelo: por cada valor ordenado de $R$ se recorren los valores de $S$ que le correspondan según $c^{52}$ (y éstos serán un conjunto ordenado de $S$ )	Costo de lectura: $b_R + b_S$ Costo de ordenación: $2 \times b \times (1 + \log_2 b)$									
	Index Join (Single Loop)	Índice sobre $S$	Recorrer $R$ accediendo por índice a $S$	$b_R + (n_R \times Z)$ Donde $Z$ : <table border="1" data-bbox="1077 1310 1508 1512"> <thead> <tr> <th>Índice</th> <th><math>Z</math></th> </tr> </thead> <tbody> <tr> <td>Secundario</td> <td><math>x + s_S</math></td> </tr> <tr> <td>Cluster</td> <td><math>x + \left\lceil \frac{s_S}{bf_S} \right\rceil</math></td> </tr> <tr> <td>Primario</td> <td><math>x + 1</math></td> </tr> <tr> <td>Hash</td> <td><math>h</math></td> </tr> </tbody> </table>	Índice	$Z$	Secundario	$x + s_S$	Cluster	$x + \left\lceil \frac{s_S}{bf_S} \right\rceil$	Primario	$x + 1$	Hash
Índice	$Z$												
Secundario	$x + s_S$												
Cluster	$x + \left\lceil \frac{s_S}{bf_S} \right\rceil$												
Primario	$x + 1$												
Hash	$h$												

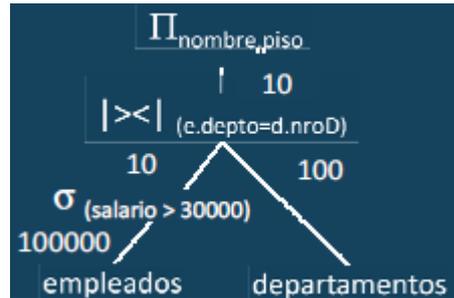
<sup>50</sup>  $x$  representa la cantidad de niveles del índice.  $s$  la cantidad de tuplas resultado de la operación

<sup>51</sup>  $M$  indica el tamaño del buffer. Es  $-2$  porque uno se usa para  $R$  y otro para  $S$ . El factor  $\left\lceil \frac{b_R}{M-2} \right\rceil$  indica la cantidad de bloques de  $S$  que se pueden leer a la vez

<sup>52</sup> Y éstos serán un conjunto ordenado de  $S$ . Por ejemplo, si  $c$  fuera de igualdad, para cada registro de  $R$  se recorren los  $S$  que le sean iguales, cuando deje de cumplirlo, es a partir de allí que debe analizarse para el siguiente registro de  $R$

- Ejemplo:

Plan Lógico



Valores del Catálogo

Parámetro	Valor
Tamaño de EMPLEADOS (tuplas)	100.000
Tamaño de DEPARTAMENTOS	100
Selectividad de $\sigma_{(\text{salario}>3000)}$	1/10.000
$Bf_{\text{empleados}}, bf_{\text{departamentos}}$	10
$bf_{\text{empleados} X \text{departamentos}}$	5
Índices sobre EMPLEADOS	B+ en salario, x=5
Todas las tablas tienen índice primario con x=1	
Asumimos que hay 3 buffers (M=3)	

Planes Físicos

	Implementación	Costo leer	Costo grabar
$\sigma_{(\text{salario}>3000)}$ Empleados	Búsqueda lineal	10.000	1
	Búsqueda binaria		No es posible
	Índice secundario	15	1
><	Loop anidado reg.	101	2
	Loop anidado bloq.	11	2
	Sort Merge	11 o 12	2
Costo total mínimo		15+11 = 26	2+1=3

El tipo de índice con el que se cuenta no garantiza datos ordenados  
 Acceder a todos los niveles del índice (5) más uno por cada tupla (hay 10 tuplas; surge de la estimación 1/10.000) =  $10 + 5$   
 $10 \times 10 + 1$   
 Es 12 si la ordenación hace una lectura más

Mejor plan

Selección con Índice Secundario y Join con Loop Anidado por Bloques, dando como resultado,  $26 + 3 = 29$  operaciones





# 7. Control de Concurrencia y Recuperación

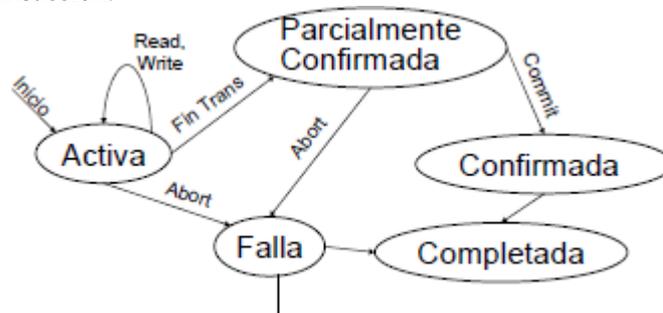
## 7.1 Transacciones e Historias

- Propiedad ACID:

Un proceso cumple con la propiedad ACID si cumple con las siguientes propiedades:

Sigla	Propiedad	Descripción
A	Atomicidad	La ejecución es considerada una unidad de trabajo
C	Consistencia	Mantiene la consistencia de la DB
I	Aislamiento	Su ejecución no interfiere con otra
D	Durabilidad	Sus modificaciones son persistentes en la DB

- Transacción:*  
Procesos concurrentes que ejecutan sobre datos compartidos y cumplen la propiedad ACID
- Estados de una Transacción:



Su resultado es como si nunca hubiese ejecutado

- Operaciones de una Transacción<sup>53</sup>:

Acción	Notación completa	Notación compacta
Leer	$read_i(X)$	$r_i(X)$
Escribir	$write_i(X)$	$w_i(X)$
Confirmar	$commit_i$	$c_i$
Abortar	$abort_i$	$a_i$

- Rollback:*  
Acción de recuperación de la DB a un estado previo a la ejecución de  $a_i$  en una Transacción
- Grafo de Seriabilidad:*  
Grafo de precedencia de las Transacciones que componen una *Historia*. Construido como:

Paso	Componente <sup>54</sup>	Acción
1	$T_i$	Nodo $T_i$
2	$w_i \dots r_j$	$T_i \rightarrow T_j$
3	$r_i \dots w_j$	$T_i \rightarrow T_j$
4	$w_i \dots w_j$	$T_i \rightarrow T_j$

<sup>53</sup>  $i$  representa el número de la Transacción y  $X$  el lugar de la DB sobre el que se opera.

La secuencia de estas operaciones, identifica una Transacción

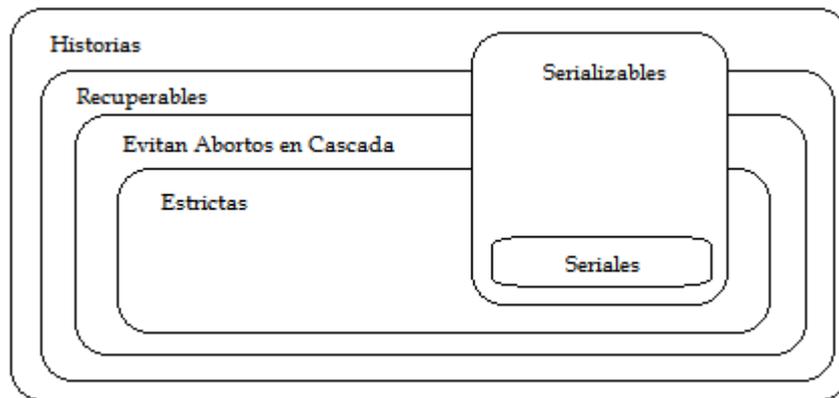
<sup>54</sup> Los  $r$  y  $w$  de esta columna tienen el mismo parámetro

- *Historia:*  
Ordenación de las operaciones de un conjunto determinado de transacciones que aparecen de forma lineal respecto de ésta<sup>55</sup>

- Clasificación de Historias:

Historia	Definición	Descripción
<i>Serializable</i>	Su Grafo de Seriabilidad es acíclico <sup>56</sup>	
<i>Recuperable</i>	Ninguna transacción confirma si no han confirmado aquellas desde las que se leyó datos previamente modificados	Los commits aparecen en el orden de flujo de datos: Se satisfacen las dependencias de confirmaciones <sup>57</sup>
<i>Evitan Abortos en Cascada</i>	Ninguna transacción lee valores escritos por otras no confirmadas	Que una transacción grabe y haga un commit antes de que otra lea
<i>Estricta</i>	Ninguna transacción lee o escribe hasta que todas las que escribieron sobre el mismo ítem fueron confirmadas	No hay un $w$ ni un $r$ si no existió sobre ese ítem modificado un $c$

- Relaciones entre las Historias:



- *Transaction Manager (TM):*  
Módulo del DBMS encargado de construir las transacciones y ordenarlas en historias serializables, recuperables y que no tengan abortos en cascada. Puede cambiar un  $c_i$  por un  $a_i$

## 7.2 Control de Concurrency

- *Lock y Unlock binarios:*  
Operación con comportamiento análogo a los semáforos de SOs cuyo dominio es una porción de la DB y su codominio la transacción en la que se ejecuta:

Operación	Notación completa	Notación compacta	Acción de bloqueo/desbloqueo sobre X
<b>Lock</b>	$lock_i(X)$	$l_i(X)$	Todas las operaciones
<b>Read Lock</b>	$read\_lock_i(X)$	$rl_i(X)$	Escrituras
<b>Write Lock</b>	$write\_lock_i(X)$	$wl_i(X)$	Todas las operaciones
<b>Unlock</b>	$unlock_i(X)$	$u_i(X)$	Todas las operaciones

<sup>55</sup> Esto es, se "intercalan" operaciones de todas las transacciones pero nunca se cambia el orden relativo dentro de la transacción al a que pertenecen

<sup>56</sup> Si posee ciclos, sufre de deadlock

<sup>57</sup> Podría pensarse como qué pasaría con transacciones que aún no confirmaron respecto de una que ya lo hizo: si el aborto de aquella impacta sobre estas, entonces, no es recuperable

- *Protocolo de Locking 2PL:*  
Reglas de locking que garanticen la seriabilidad:

Variante	Descripción	Características	Ejemplo
<b>Básico</b>	A lo largo de una Historia, primero están todos los $l$ (fase de expansión) y luego y luego todos los $u$ (fase de contracción)	$\uparrow$ Fácil $\downarrow$ Deadlock <sup>58</sup>	$\underbrace{rl_i(X) \dots wl_j(Y)}_{\text{Expansión}} \dots \underbrace{u_i(X) \dots u_j(Y)}_{\text{Contracción}}$
<b>Conservador</b>	Declarar primero todos los $l$	$\uparrow$ No hay Deadlock $\downarrow$ Exige declarar todo lo que se utilizará	$rl_i(X)wl_j(Y) \dots u_i(X)u_j(Y)$
<b>Estricto</b>	No se liberan $wl$ hasta confirmar la $T$	$\uparrow$ Historias <i>Estrictas</i> $\downarrow$ Deadlock	$\dots wl_i(X) \dots c_i u_i(X)$
<b>Riguroso</b>	No se liberan ni $wl$ ni $rl$ hasta confirmar la $T$	$\uparrow$ Fácil (respecto del Estricto) $\downarrow$ Deadlock	$\dots rl_i(X)wl_j(Y) \dots c_i u_i(X)u_j(Y)$

- *TimeStamp (TS):*  
Para cada ítem  $X$  se define un  $Read\_TS(X)$  y un  $Write\_TS(X)$  que informa el TS de la transacción más joven que leyó y escribió sobre él respectivamente

- Control de Concurrencia basado en TS:

Acción que pretende $T$	Condición para su ejecución	$TS(T)$	Acción
$w(X)$	Nadie leyó o escribió posteriormente a $T$	$< Read\_TS(X)$	$T$ es reiniciado con un TS actualizado
		$< Write\_TS(X)$	
$r(X)$	Nadie escribió posteriormente a $T$	$< Write\_TS(X)$	$T$ es reiniciado con un TS actualizado
		Caso contrario	

- *Multiversión:*  
Se mantienen varias versiones de cada ítem y se selecciona la adecuada para cada transacción:

$$X = \{ ( X_i, Read_{TS(X_i)}, Write_{TS(X_i)} ) \}_{i=1, \dots, n}$$

- Control de Concurrencia basado en Multiversión:

Acción que pretende $T$	Condición para su ejecución	$TS(T)$	Acción
$w(X)$	Nadie leyó posteriormente a $T$	$Write\_TS(X_i) \leq TS(T) < Read\_TS(X_i)$	$T$ es reiniciado con un TS actualizado
		Caso contrario	Se crea $(X_j, TS(T), TS(T))$
$r(X)$	Ninguna: Se utiliza la última escritura más cercana a $T$	$< Write\_TS(X)$	$T$ es reiniciado con un TS actualizado
		Caso contrario	$Read\_TS(X) = \max\{TS(T), Read\_TS(X)\}$

### 7.3 Otros aspectos de Concurrencia

- *Granularidad:*  
Refiere al detalle con el que se define un ítem  $X$ : tupla, registro, tabla, bloque de disco...
- Granularidad y performance:

Granularidad	Impacto
<b>Mayor</b>	Menor concurrencia
<b>Menor</b>	Mayor overhead

<sup>58</sup> En esta columna, esto significa "Susceptible a..."

- *Registros Fantasmas:*  
Es un registro que, a la espera de un desbloqueo, no es insertado pese a que es de interés para la transacción bloqueante
- Detección de Deadlocks:  
Mantener un grafo de espera que detecte la existencia de deadlocks
- Solución a Deadlocks:

**Timeout** Si una transacción espera por mucho tiempo, reiniciarla

Basados en TS	Estrategia	Descripción
	<b>Wait-Die</b>	La transacción más nueva del deadlock debe abortar y recomenzar con el mismo TS
	<b>Wound-Wait</b>	La transacción más nueva del deadlock debe esperar por la más vieja

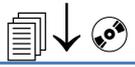
- Starvation:  
Wait-Die y Wound-Wait evitan el problema que, generalmente, pueden originarse por malos mecanismos de "selección de víctima" (puede suceder por su reinicio indefinido)

## 7.4 Recuperación

- *Log:*  
Registro de actividad sobre la DB
- Administración de logs:  
Deben respaldarse frecuentemente y almacenarse en lugares distintos a la DB
- Valores posibles de un ítem:

Sigla	Valor	Valor del ítem respecto de la actualización	Log de <sup>59</sup>
<b>BFIM</b>	BeFore IMage	Antes	Undo
<b>AFIM</b>	AFTer IMage	Después	Redo

- Ciclo de vida de una instrucción de una transacción:

Paso	Descripción	Esquema
1	Se hacen los cambios en un buffer de memoria	
2	Se registran los cambios en el log	
3	La transacción confirma	
4	Se baja el log al disco	
5	Se bajan los cambios al disco	

- Técnicas de recuperación:

Actualización	Operativa	Ante un abort
<b>Diferida</b>	Cada transacción trabaja en una sandbox que es enviada a disco luego del commit	No debe realizarse acción alguna
<b>Inmediata</b>	La base es actualizada antes del commit	Se deben deshacer las operaciones de la transacción valiéndose del log
<b>Shadow Paging<sup>60</sup></b>	Cada transacción mantiene en memoria un conjunto de punteros hacia los valores originales y los va cambiando a nuevas páginas de memoria a medida que surjan modificaciones	Recupera el conjunto de punteros original

- Write-Ahead Loggin (WAL):*  
Protocolo que garantiza que el log con el BFIM está en el disco antes de grabar la operación en la DB
- Checkpoint:*  
Registro del log que indica que todos los buffers modificados de la DB están en disco:

Paso	Descripción	Esquema
1	Suspender transacciones en curso	
2	Bajar buffers modificados	↓ ⚡
3	Registrar checkpoint	✓ 📄
4	Bajar log	↓ 📄
5	Reanudar transacciones	▶

- Rollback:*  
Se debe realizar cuando una transacción aborta. Si otras transacciones se basaron en ella para continuar, éstas también deben ser abortadas, provocando un aborto en cascada

- Algoritmos de recuperación<sup>61</sup>:

Algoritmo	Acciones de <i>T</i> confirmada	Acciones de <i>T</i> no confirmada
<b>Basado en Actualización Diferida</b>	Se reconstruye el resultado utilizando el log	Se relanza al finalizar la recuperación
<b>Basada en Recuperación Inmediata</b>	Se reconstruye el resultado utilizando el log	Se revierten los cambios empleando los valores anteriores en el log. Se relanza la transacción
<b>Shadow Paging</b>	Se graban los punteros modificados	Se graban los punteros originales

<sup>61</sup> Todos ellos asumen la generación de historias estrictas y serializables