

# Arquitectura de un sistema de base de datos

Clase 01

IIC 3413

Prof. Cristian Riveros

# Outline

Modelo Relacional

SQL

Algebra Relacional

Arquitectura

# Outline

**Modelo Relacional**

SQL

Algebra Relacional

Arquitectura

# Modelo relacional

## Definición

Una **relación**  $R$  esta compuesto por:

1. Esquema:  $\text{name}(\text{att}_1 : D_1, \dots, \text{att}_n : D_n)$ 
  - $\text{name} :=$  nombre de la relación  $R$ .
  - $\text{att}_i :=$  nombre del atributo  $i$ .
  - $D_i :=$  dominio del atributo  $i$ .
  - $n :=$  aridad o cantidad de atributos de  $R$ .
2. Instancia: subconjunto finito de  $D_1 \times \dots \times D_n$ .

Dada una relación  $R$ , usualmente denotaremos

- $R$ : como el nombre y instancia de  $R$  (si no hay confusión).
- $\text{schema}(R)$ : como el esquema de  $R$ .
- $\text{arity}(R)$ : como la aridad de  $R$ .
- $\text{dom}(\text{att}_i)$ : como el dominio  $D_i$  del atributo  $\text{att}_i$ .

# Propiedades de una relación

El dominio de un atributo puede ser cualquier tipo primitivo:

- Números
- Strings
- Fechas, ...

Una relación  $R$  se ve como una tabla con filas y columnas:

<b>att<sub>1</sub></b>	<b>att<sub>2</sub></b>	...	<b>att<sub>n</sub></b>
$v_1$	$v_2$	...	$v_n$
$v'_1$	$v'_2$	...	$v'_n$
$\vdots$	$\vdots$		$\vdots$

donde  $(v_1, \dots, v_n)$ ,  $(v'_1, \dots, v'_n)$ , ... son tuplas de  $R$ , PERO...

- El orden de las filas no es importante.
- Todas las filas son distintas.

# Base de Datos (BD) relacional

## Definición

1. Una **BD relacional**  $\mathcal{D}$  es un conjunto finito de relaciones  $R_1, \dots, R_m$  cada una con un nombre distinto.

$$\mathcal{D} = \{R_1, R_2, \dots, R_m\}$$

2. Un **esquema relacional**  $\mathcal{S}$  de una BD  $\mathcal{D}$  es el conjunto de esquemas de las relaciones  $R_1, \dots, R_m$ .

$$\text{schema}(\mathcal{D}) = \{\text{schema}(R_1), \text{schema}(R_2), \dots, \text{schema}(R_m)\}$$

Otras definiciones importantes:

- $\text{dom}(\mathcal{S})$  = todas las BD que tienen como esquema a  $\mathcal{S}$ .
- $|R|$  = numero de tuplas en  $R$  (o cardinalidad).

# Base de Datos (BD) relacional

## Ejemplo

### ■ Esquema:

Players(*pName* : CHAR, *pTeam* : CHAR, *pBirth* : DATE, *pPosition* : INT)

Teams(*tName* : CHAR, *tFans* : INT, *tBirth* : DATE)

### ■ Instancia:

<b>pName</b>	<b>pTeam</b>	<b>pBirth</b>	<b>pPosition</b>
Alexis Sanchez	Inter de Milan	Dec 19, 1988	7
Gary Medel	Bologna	Ago 3, 1987	5
		<b>tName</b>	<b>tFans</b>
		<b>tBirth</b>	
		Bologna	38M
		Inter de Milan	80M

# Restricciones de Integridad

## Definición

Una **restricción de integridad** es una condición  $\varphi$  que restringen los datos que pueden ser almacenados en una BD relacional.

Una BD  $\mathcal{D}$  es **valida** con respecto a  $\varphi$  si satisface las restricción  $\varphi$ .

$$\mathcal{D} \models \varphi$$

## Ejemplos

- Restricciones de dominio.
- Llaves (Keys).
  - Primary Key.
  - Foreign Key.
- Condiciones generales.

Solo se permiten BD validas según las restricciones de integridad.



# Extracción de datos basado en consultas

**Consulta:** función  $f : \text{dom}(S) \rightarrow D$

- $S$  es una esquema relacional,
- $D$  es un dominio cualquiera ( $\mathbb{N}$ , relaciones, etc).

**Lenguaje de consultas:** conjunto de expresiones **sintácticas** que definen consultas por medio de una **semántica**.

Dos lenguajes de consultas importantes para nosotros:

1. SQL
2. Algebra relacional

# Outline

Modelo Relacional

**SQL**

Algebra Relacional

Arquitectura

# SQL (Structured Query Language)

- Standard mundial para consultas BD relacional.
- Lenguaje de consultas declarativo.
- Basado en calculo relacional (lógica de primer orden).
- Múltiples componentes del lenguaje:
  - Data definition language (DDL).  
Ejemplo: CREATE, ALTER, etc...
  - Data manipulation language (DML).  
Ejemplo: SELECT, INSERT, DELETE, etc...
  - Data control language (DCL).  
Ejemplo: GRANT, REVOKE, etc...

# Consultas SQL

Forma basica:

SELECT	< atributos >
FROM	< relaciones >
WHERE	< condiciones >

< atributos >: lista de atributos

- atributo
- Relacion . atributo
- Relacion . atributo AS nuevo\_nombre

< relaciones >: lista de relaciones

< condiciones >: lista de condiciones booleanas

- $\varphi_1$  AND  $\varphi_2$ ,  $\varphi_1$  OR  $\varphi_2$ , NOT  $\varphi$
- atributo<sub>1</sub> ~ atributo<sub>2</sub> donde  $\sim \in \{=, \leq, \geq, \dots\}$ .
- atributo<sub>1</sub> ~ *constante* donde  $\sim \in \{=, \leq, \geq, \dots\}$ , etc...

# Ejemplos de consultas SQL

## Ejemplos

**Players (P):**

pId	pName	pYear	pGoals
1	x	1987	100
2	y	1990	59
3	y	1985	88
NULL	z	1983	110

**Matches (M):**

mId	mStadium	mGoals
1	Nacional	3
2	Nacional	2
2	San Carlos	3
5	Monum.	4

- ```
SELECT  pName, pYear
FROM    Players
WHERE   pGoals ≥ 100
```
- ```
SELECT  P.pName AS Jugador, P.pYear AS Año
FROM    Players AS P
WHERE   pYear = 1990
```

# Ejemplos de consultas SQL

## Ejemplos

**Players (P):**

pId	pName	pYear	pGoals
1	x	1987	100
2	y	1990	59
3	y	1985	88
NULL	z	1983	110

**Matches (M):**

mId	mStadium	mGoals
1	Nacional	3
2	Nacional	2
2	San Carlos	3
5	Monum.	4

- ```
SELECT  DISTINCT pName AS Nombre, mStadium AS Estadio
FROM    Players, Matches
WHERE   pId = mId
```
- ```
SELECT  pName, AVG(mGoals)
FROM    Players, Matches
WHERE   pId = mId
GROUP BY pId, pName
```

# Consultas anidadas y/o correlacionadas

**Consultas anidada:** consulta que contienen una subconsulta embebida en:

- WHERE
- FROM
- HAVING

## Ejemplos

- ```
SELECT  pName, pGoals
FROM    Players
WHERE   pGoals = ( SELECT  MAX(pGoals)
                   FROM    Players )
```
- ```
SELECT  DISTINCT pName, pYear
FROM    Players, ( SELECT  mld
                   FROM    Matches
                   WHERE   mGoals ≥ 3 )
WHERE   pld = mld
```

# Consultas anidadas y/o correlacionadas

**Consulta correlacionada:** consulta anidada que contienen una subconsulta embebida con una referencia a la consulta externa.

## Ejemplo

```
SELECT  pName, pYear
FROM    Players AS P
WHERE   1 ≤ ( SELECT  AVG(mGoals)
              FROM    Matches AS M
              WHERE    M.mld = P.pld )
```

¿cuál de las consultas anteriores  
puede ser definida con una consulta NO anidada?



# Outline

Modelo Relacional

SQL

**Algebra Relacional**

Arquitectura

# Algebra relacional

- Lenguaje de consultas procedural.
- Basado en operadores relacionales (algebra).

¿cuál es la utilidad del algebra relacional en BD?

Ventajas:

- Fácil de componer.
- Fácil de optimizar.

# Operadores relacionales

Selección:  $\sigma_{\theta}(R)$ .

- $\theta$  es una combinación booleana ( $\wedge, \vee$ ) de terminos:

$atributo_1$  op  $atributo_2$   
 $atributo$  op  $constante$

con  $op \in \{=, \leq, \geq, <, >\}$ .

Proyección:  $\pi_L(R)$ .

- $L$  = lista de atributos.

Operaciones de conjunto:

- Union:  $R_1 \cup R_2$
- Intersección:  $R_1 \cap R_2$
- Diferencia:  $R_1 - R_2$

# Operadores relacionales

Joins:

- Producto cruz:  $R_1 \times R_2$
- $\theta$ -join:  $R_1 \bowtie_{\theta} R_2 = \sigma_{\theta}(R_1 \times R_2)$ 
  - $\theta$  es una combinación booleana ( $\wedge, \vee$ ) de terminos:

$$\begin{array}{lcl} \textit{atributo}_1 & \text{op} & \textit{atributo}_2 \\ \textit{atributo} & \text{op} & \textit{constante} \end{array}$$

con  $\text{op} \in \{=, \leq, \geq, <, >\}$ .

- Equi-join:  $R_1 \bowtie_{\phi} R_2 = \sigma_{\phi}(R_1 \times R_2)$ 
  - $\phi$  solo contine igualdades.
- Natural-join:  $R_1 \bowtie R_2 = \sigma_{\phi}(R_1 \times R_2)$ 
  - $\phi = \bigwedge_{a \in \text{att}(R_1) \cap \text{att}(R_2)} R_1.a = R_2.a.$

# Ejemplo de consultas en algebra relacional

## Ejemplos

**Players (P):**

id	name	year	goals
1	x	1987	100
2	y	1990	59
3	y	1990	88
4	z	1983	110

**Matches (M):**

id	stadium	mgoals
1	Nacional	3
2	Nacional	2
2	San Carlos	3
NULL	Monum.	4

- $\pi_{\text{name,year}}(P)$
- $\sigma_{\text{goals} \geq 100 \vee \text{goals} \leq 60}(P)$
- $\pi_{\text{goals}}(\sigma_{\text{year}=1990}(P))$
- $P \bowtie_{P.id=M.id} M$
- $P \bowtie M$

# Semántica algebra relacional

**Players (P):**

id	name	year	goals
1	x	1987	100
2	y	1990	59
3	y	1990	88
4	z	1983	110

**Matches (M):**

id	stadium	mgoals
1	Nacional	3
2	Nacional	2
2	San Carlos	3
NULL	Monum.	4

¿cuál es el resultado de la consulta  $\pi_{\text{name, year}}(P)$ ?

Dos tipos de semánticas para algebra relacional:

- Set semantics.
- Bag semantics (SQL).

¿por qué usar set semantics o bag semantics?

# Algebra relacional extendida

Rename:  $\rho_{old\_att \rightarrow new\_att}(R)$

Eliminación de duplicados:  $\delta(R)$

Group-by con agregación:  $\gamma_{G,A}(R)$

- G: lista de atributos ha agrupar.
- A: lista de elementos de la forma:

$$f(agg\_att) \rightarrow new\_att$$

con  $f \in \{MIN, MAX, SUM, AVG, \dots\}$ .

Sorting (ordenar):  $\tau(R)$

# Algebra relacional extendida

## Ejemplos

**Players (P):**

id	name	year	goals
1	x	1987	100
2	y	1990	59
3	y	1990	88
4	z	1983	110

**Matches (M):**

id	stadium	mgoals
1	Nacional	3
2	Nacional	2
2	San Carlos	3
NULL	Monum.	4

- $P \bowtie \rho_{mgoals \rightarrow goals}(M)$
- $\delta(\pi_{name, year}(P))$
- $\gamma_{year, AVG(goals) \rightarrow GperY}(P)$



¿cómo evaluarían esta consulta Q?

Players(pld, pName, pYear)

Matches(mld, mStadium, mDate)

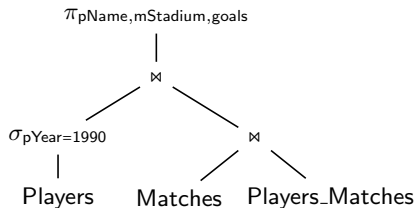
Players\_Matches(pld, mld, goals)

$\pi_{pName, mStadium, goals} ($

$\sigma_{pYear=1990} (Players) \bowtie ($

$(Matches \bowtie Players\_Matches)))$

**Plan lógico** (= árbol de parsing) de Q:



El **plan lógico** será nuestro punto inicial para la **evaluación** de la consulta.

¿cómo convertirían esta consulta a algebra relacional?

```
SELECT  pName, mStadium, goals
FROM    Players AS P, Matches AS M, Players_Matches AS PM
WHERE   P.pId = PM.pId AND PM.mId = M.mId AND
        P.pYear  $\geq$  1985
```

# Desde SQL a Algebra relacional

## Idea (intuitiva)

```
SELECT  att1, ..., attn
FROM    R1, ..., Rm
WHERE    $\varphi$ 
```

1. Combine las relaciones en el FROM con **productos cruz**:

$$R_1 \times \dots \times R_n$$

2. Aplique **selección** con la condición dada en el WHERE:

$$\sigma_{\varphi}(R_1 \times \dots \times R_n)$$

3. Aplique **proyección** con los atributos dados en el SELECT:

$$\pi_{att_1, \dots, att_n}(\sigma_{\varphi}(R_1 \times \dots \times R_n))$$

¿es este resultado un “buen” plan de evaluación de la consulta?

# Outline

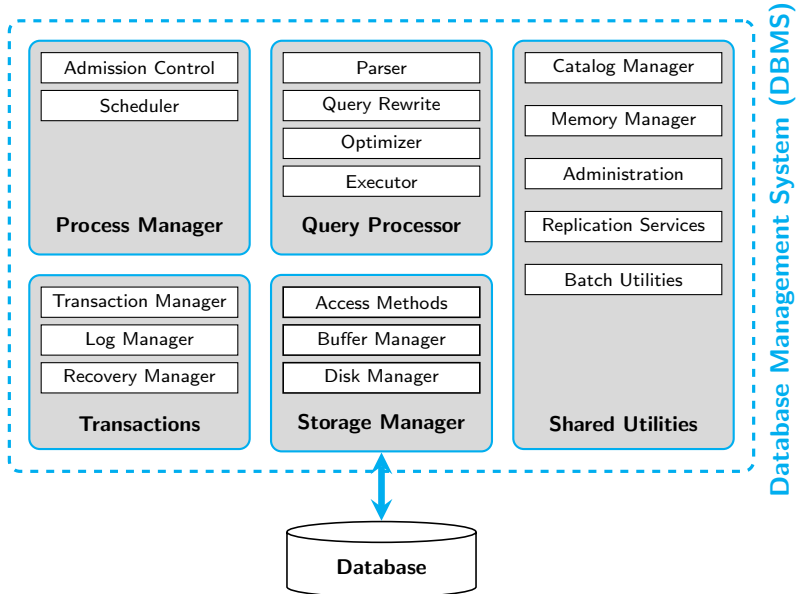
Modelo Relacional

SQL

Algebra Relacional

**Arquitectura**

# Arquitectura de un Sistema de Bases de Datos (DBMS)

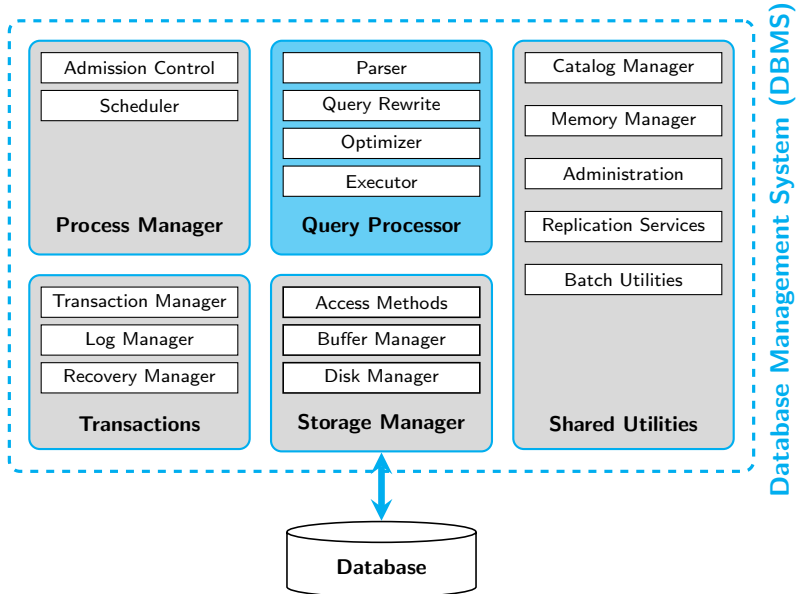


# La vida de una consulta SQL

Considere la siguiente consulta:

```
Q  =  SELECT  pName, mStadium, goals
      FROM    Players AS P, Matches AS M, Players_Matches AS PM
      WHERE   P.pId = PM.pId AND PM.mId = M.mId AND
              P.pYear  $\geq$  1985
```

# Procesador de consultas



# Procesador de consultas

## Paso 1: Parser

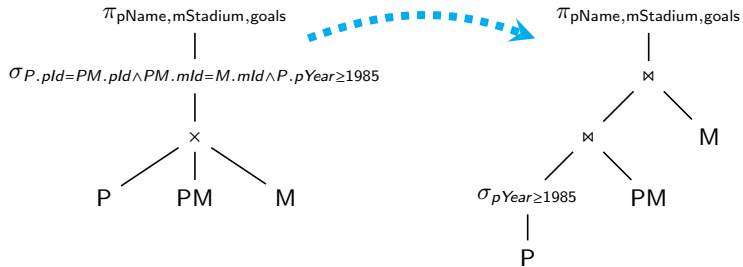
- Valida la consulta (correctitud, autorización, etc).
- Convierte consulta en formato interno (álgebra relacional).
- Optimizaciones menores (numéricas, etc).

## Paso 2: Reescritura de consulta

- “Desanidación” de la consulta (flattening).
- Reescribe consulta aplicando reglas de álgebra relacional.
- Crea un set de planes lógicos para ser optimizados.



# Procesador de consultas



# Procesador de consultas

## Paso 3: Optimizador

- Encuentra el **plan físico** mas eficiente para ejecutar la consulta.
- El optimizador debe considerar:
  - Tamaño de cada relación y distribución de sus datos.
  - Distintos accesos a las relaciones (índices, etc).
  - Distintos planes lógicos para la misma consulta.
  - Distintos algoritmos para un mismo operador relacional.

## Paso 4: Ejecución

- Ejecuta el plan óptimo encontrado por Optimizador (**pipelining**).

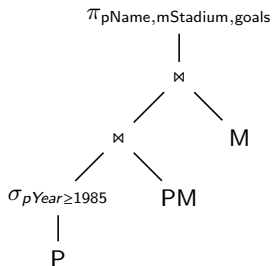
# Plan físico

El **Plan físico** es parecido al **Plan lógico** pero:

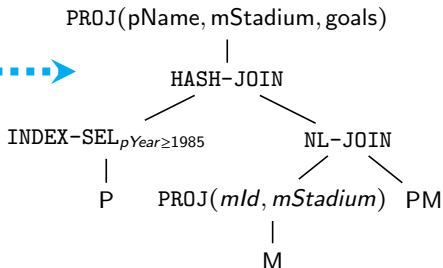
- Selecciona el **access path** para cada relación
  - Uso de índices o (simplemente) file scan.
- Decide el algoritmo ha utilizar por cada operador.
- Planifica secuencia y orden de los operadores.

# Plan lógico vs. Plan físico

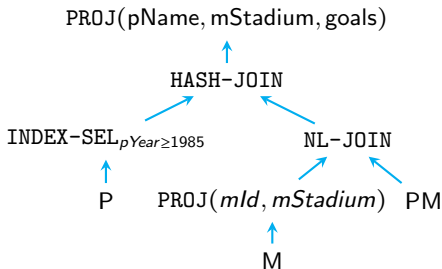
**Plan lógico**



**Plan físico**

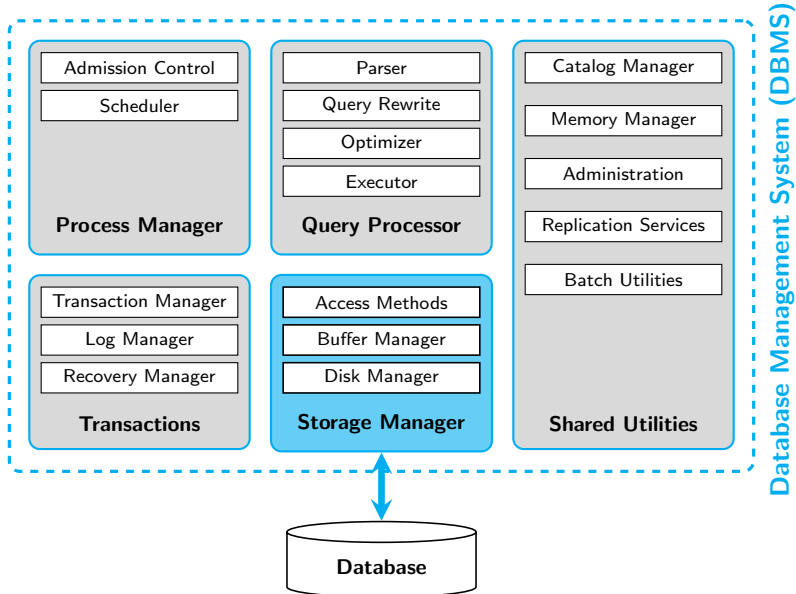


# Pipelining (Ejecución)



- Ejecución en serie del plan físico.
- Cada operador retorna tuplas a su operador **padre**.
- Cada tupla es retornada “as soon as possible”.

# Manejador de almacenamiento



# Manejador de almacenamiento

## Disk Manager

- Acceso para almacenamiento secundario (disco duro, etc).
- Organizador de tuplas (**Heapfiles**).

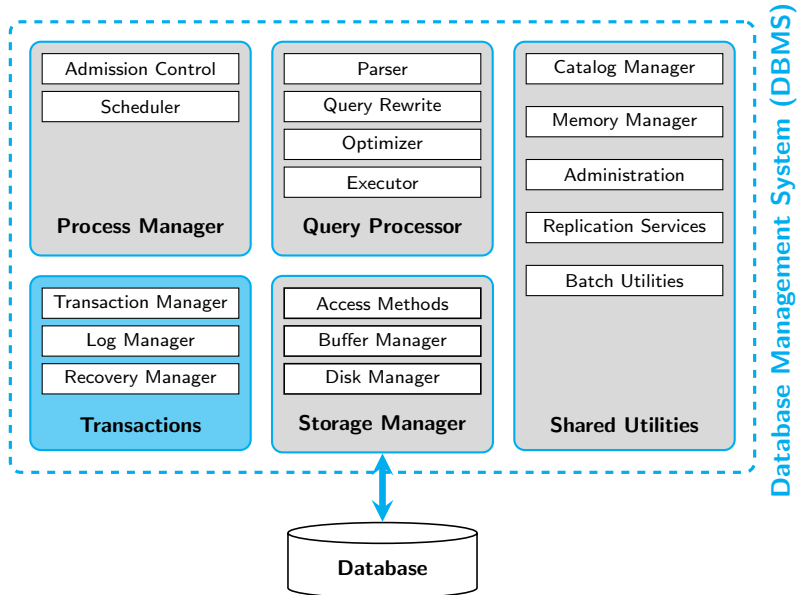
## Buffer Manager

- Manejo de datos en memoria ( $|memoria| \ll |disco\ duro|$ ).
- Optimiza la cantidad de accesos al disco duro (I/O).
- Importante para transacciones (ACID).

## Access Methods

- Todo tipo de índices.
- Organización eficiente de los datos.

# Transacciones





# Transacciones

**ACID** = **A**tomicity  
**C**onsistency  
**I**solation  
**D**urability

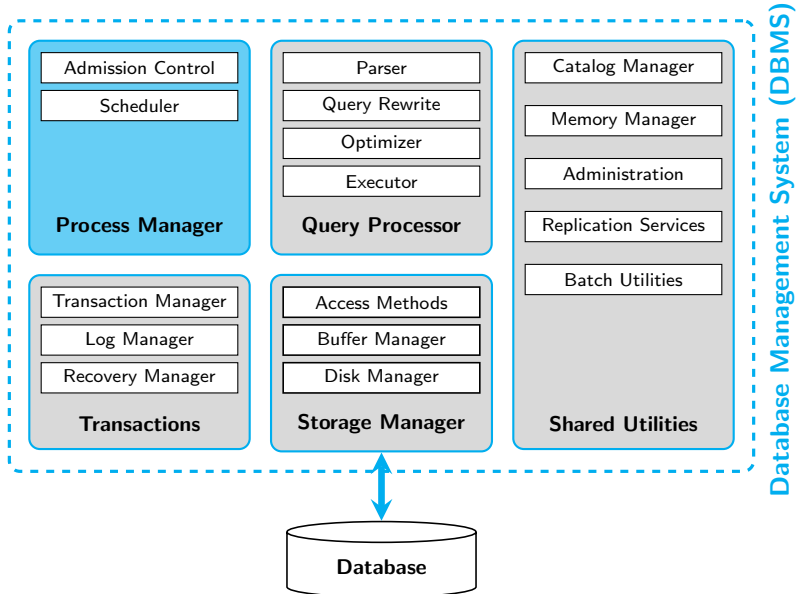
## Transaction Manager

- Encargado de asegurar **I**solation y **C**onsistency.

## Log y Recovery Manager

- Encargado de asegurar **A**tomicity y **D**urability.

# Manejador de procesos



# Manejador de procesos

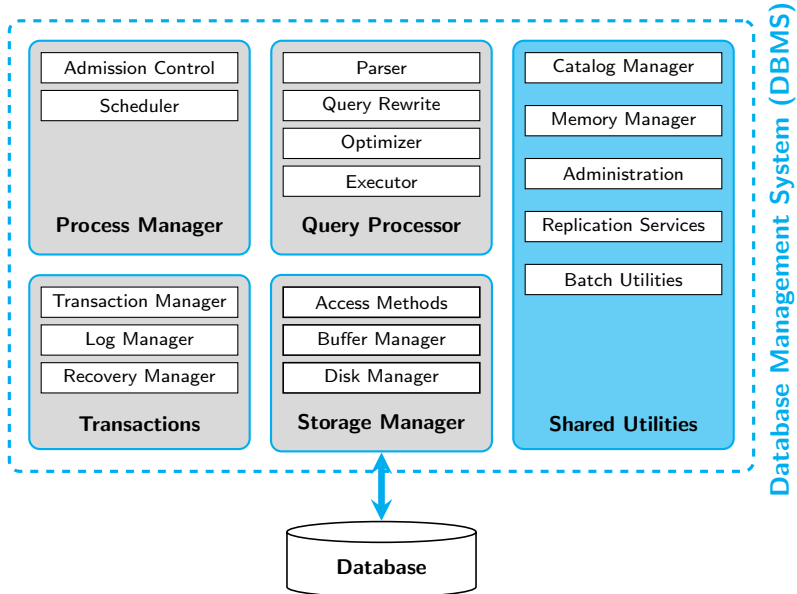
## Control de Admisión

- Manejo de usuarios (autenticación, permisos, etc).
- Encargado de asegurar la cantidad de recursos (trashing).

## Scheduler

- Manejo de procesos y threads.
- Paralelización de consultas.

# Utilidades



# Varias utilidades (alguna de ellas)

## Catalogo

- Esquema y metadata de los datos.

## Manejador de memoria

- Asignador de memoria para otros componentes.

# Próxima clase empezaremos desde el **Storage Manager**

