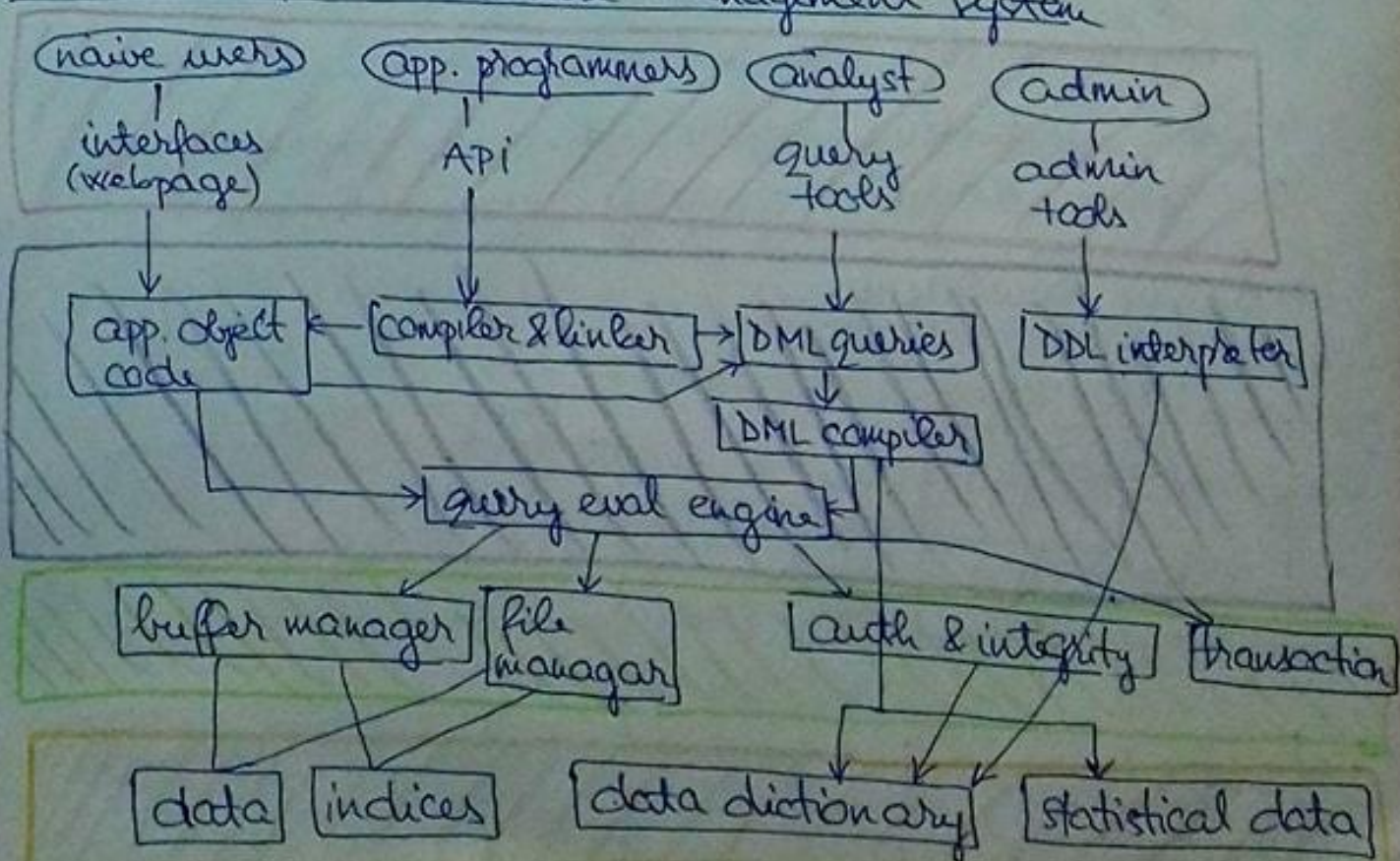# Database systems
## FÉLÉVES ÖSSZEFOGLALÓ VIZSGÁRA

1. <u>Introduction, course outline</u>
   - What is a DBMS ? ✓
   - What is data modelling ? ✓
   - Abstraction levels of a database ✓
   - Creating databases (abstraction level, logical struct., DDL) ✓
   - Getting data from database : QL → RELATIONAL ALGEBRA ✓
   - Manipulating data in database : DML → SQL ✓
   - Storing data : authorization, recovery → UNDO/REDO LOGGING SHADOWING
   - Multiuser environment : transaction management, concurrency control
   - Developement lifecycle ✓
       ↓
       LOCKING MECHANISMS
       TIMESTAMPING
   - Active elements : constraints, triggers
   - Database architectures ✓
   - Functional dependencies and normalisation — ✓

2. <u>Structure of a database management system</u>



-1-

## What is a database:

- Structured data in relational storage (properties)
- serves multiple users:
  - access, insert, modify data → QL, DML → authority management
  - simultaneous access → transaction manager
- history (recovery)
- File Management (physical layer)

## Data modelling, abstraction levels

WORLD → abstraction : no values, just relations
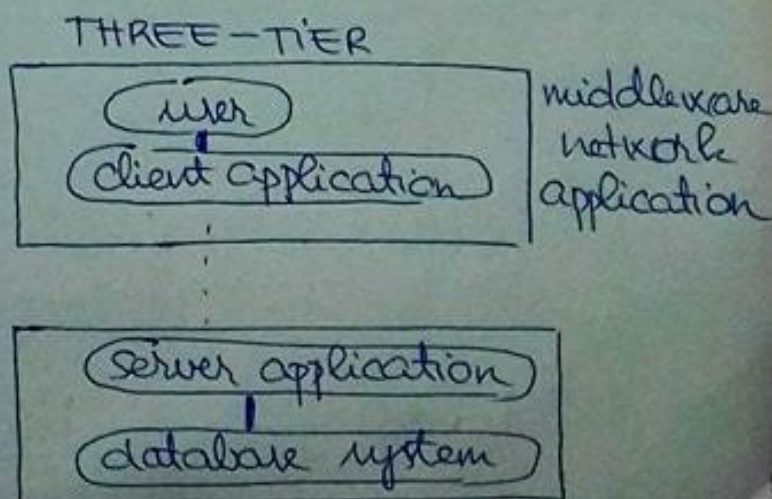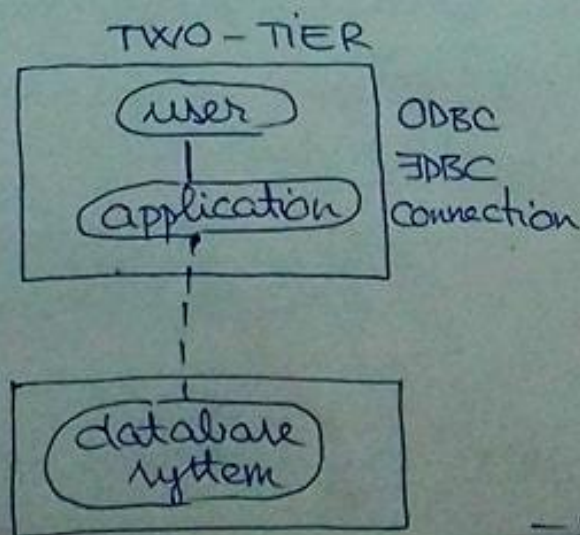database INSTANCE → satisfying model

Levels:
- view level : accessible by naive users and apps
- logical/conceptual level : DB admins
- physical level : DBMS programmers :P

## Creating databases/developement lifecycle

1. specification (~~high level design~~) : requirements, constraints
2. high level design : entity-relationship model
3. implementation design : relational model
4. low level design : files, indices → ~~implementation~~
5. system model : complete documentation → implement.

## Database architectures

TWO-TIER



THREE-TIER

# Creating databases, modeling data (E-R)
## HIGH LEVEL DESIGN

**Entity**: object distinguishable from others in the microworld, through its attributes

**Entity set**: group of similar objects, sharing the same properties called attributes

**Relationship**: logical connection between two or more entity sets (association among entities)

**Relationship (set)**: mathematical relation:

$$R = \{(e_1, e_2, \ldots, e_N) \mid e_1 \in E_1, e_2 \in E_2, \ldots, e_N \in E_N\}$$

where $(e_1, e_2, \ldots, e_N)$ is a relationship (n-ary tuple)

**Relation** $r(R) \subseteq E_1 \times E_2 \times \ldots \times E_N \rightarrow$ table

**Relation schema**: $R(E_1, E_2, \ldots, E_N) \rightarrow$ table header

**Types of relationship** (based on degree — number of elements)

- one to one relationship
- one to many relationship
- many to many relationship
  - total participation
  - partial participation

**Keys**:

SUPERKEY: set of attributes uniquely identifying the entity

CANDIDATE KEY: minimal set of attributes, forming an SK

PRIMARY KEY: chosen CK.

> FOREIGN KEY:
> - PK of another entity set
> - must exist in the other (constraint)

**Weak entity set**: has not enough attributes for identification

- → joined by so-called weak one to many r. to a strong set.
- → discriminator: unique in context.

**Generalization, specialization** (ISA): the sub entity inherits all properties of its anchestor

- may or may not be a partition of the superset

<u>Domain constraint</u> : set of ~~all~~ values an attribute can take

<u>NULL values and three-valued logic</u>

NULL value : data is not known, not available (permission contr.)
    or has no meaning in the given context.
Arithmetic interpretation : operations are undefined

Logical operations :
    TRUE — 1          AND : min(...)
    NULL — 1/2        OR  : max(...)
    FALSE — 0

<u>Relational model</u> : IMPLEMENTATION DESIGN

— transformation of the high-level design
— entity sets ⇒ tables, attributes ⇒ columns
— binary relationships ⇒ ┌─ One-to-one : foreign key on
                                any of the sides
                                one-to-many : foreign key on
                                the many side
                                many-to-many : separate
                                table with two foreign keys
— non-binary relationship ⇒ artificial (weak) entity set
    and multiple binary relations
— ISA relationship types :
    — condition-defined vs. user-defined
    — disjoint vs. overlapping
    — total or partial
    — aggregated

# Relational algebra. Getting data from database

Relation schema: $R(A_1, A_2, ..., A_n)$ where $A_n$ are attributes

Relation: $r(R) \subseteq D_1 \times D_2 \times ... \times D_n$

Relation instance: current values of a relation (table)

Keys (revision):

SK: K is SK of R if values of K identify a unique
    tuple of each possible relation in $r(R)$

CK: $\forall$ subset of CK does not identify + is an SK

Strong, weak and relationship entity set.

1. SELECT: $\sigma_p(r) = \{t \mid t \in r \text{ and } p(t) \text{ is true}\}$

    p: Selection predicate, propositional calculus

2. PROJECT: $\pi_{A_1, A_2, ..., A_k}(r)$: deleted columns,
    removed duplicates

3. UNION: $r \cup s = \{t \mid t \in s \text{ or } t \in r\}$

    constraint: r, s compatible ⟵ same arity
                        convertible data type

4. DIFFERENCE: $r - s = \{t \mid t \in r \text{ and } t \notin s\}$

    constraint: r, s compatible

5. CARTESIAN: $r \times s = \{(t, q) \mid t \in r \text{ and } q \in s\}$

    constraint: disjoint attributes

6. RENAME: $\rho_{X(B_1, B_2, ..., B_n)}(E) \longrightarrow$ returns $E(A_1, ..., A_n)$
    under the name/schema $X(B_1, B_2, ...)$

## ADDITIONAL OPERATORS

7. SET INTERSECT: $r \cap s = r - (r - s)$

## 8. NATURAL JOIN

$$r \bowtie s = \pi_{r.A_1, r.A_2, \ldots, r.A_n, s.B_{n+1}, s.B_{n+2}, \ldots, s.B_k} \left( \sigma_{\ldots} (r \times s) \right)$$

## 9. DIVISION

$$R = (A_1, \ldots, A_n, B_1, \ldots, B_m)$$
$$S = (B_1, \ldots, B_m)$$

$$r \div s = \{t \mid t \in r \text{ and } \forall u \in s \ tu \in r\}$$
$$r \div s = \pi_{R-S}(r) - \pi_{R-S}\left((\pi_{R-S}(r) \times s) - r\right)$$

## 10. ASSIGNMENT : convenience, nothing else

## EXTENDED OPERATIONS

## 11. GENERALIZED PROJECTION : operations on attributes are also permitted (arithmetical, string)

## 12. AGGREGATE FUNCTIONS

$$(G_1 G_2 \ldots G_n \ \mathcal{g} \ F_1(A_1), F_2(A_2), \ldots, F_n(A_n)) (E)$$

$G_1, G_2, \ldots, G_n$ : attributes on which to group
$F_1, F_2, \ldots, F_n$ : operations (avg, min, max, sum, count)
$A_1, A_2, \ldots, A_n$ : attributes

Convenience : renaming with "as" of $F(A_n)$ is permitted.

## 13. OUTER JOIN (~~inner~~, left, right, full)
  - nonmatching ~~relatio~~ tuples filled with **null values**

## Manipulating data with relational algebra

1. DELETE :   $r \leftarrow r - E$

2. INSERT :   $r \leftarrow r \cup E$

3. UPDATE :   $r \leftarrow \pi_{F_1, F_2, \ldots, F_l}(r)$

# Functional dependencies

## Outline

- functional dependency is a _semantic_ concept, used ↗ defined by the world, being modeled
  to define formal measures on the goodness of a relational design
- can be defined as a constraint on legal relation instances
- used to define normal forms

## Formal definition

"Tuple function" : $t : \{A_1, A_2, ..., A_n\} \rightarrow D_{A_1} \cup D_{A_2} \cup ... \cup D_{A_n}$

    ↳ example : $t[ID] = $ "Bus123", ← value in the tuple

Functional dependency :

$\alpha = \{A\ell_1, A\ell_2, ..., A\ell_n\}$    $\beta = \{A\ell_1, A\ell_2, ..., A\ell_n\}$

$\alpha \rightarrow \beta$ means that $t_1[\alpha] = t_2[\alpha] \Longrightarrow t_1[\beta] = t_2[\beta]$

Trivial dependency : satisfied by all tuples.

Partial dependency : attributes depend on part of the key, not all attributes of it. [this should be avoided]

Armstrong axioms

3. Transitivity    : $\alpha \rightarrow \beta$ and $\beta \rightarrow \gamma \Rightarrow \alpha \rightarrow \gamma$

2. Augmentation : $\alpha \rightarrow \beta \Rightarrow \alpha\gamma \rightarrow \beta$

1. Reflexivity    : $\alpha \rightarrow \beta$ if $\beta \subseteq \alpha$    Proof : trivial

Rules derived from this :

4. Pseudotransitivity : $\alpha \rightarrow \beta$ and $\beta\gamma \rightarrow \delta \Rightarrow \alpha\gamma \rightarrow \delta$

5. Union    : $\alpha \rightarrow \beta$ and $\alpha \rightarrow \gamma \Rightarrow \alpha \rightarrow \beta\gamma$

5. Decomposition : $\alpha \rightarrow \beta\gamma \Rightarrow \alpha \rightarrow \beta$ and $\alpha \rightarrow \gamma$

## [Derivation] (bizonyítás, levezetés, bővítés — nem $\frac{d}{dx}$ :P)

**Implication**: $FD_1 \models FD_2$ if every instance legal under $FD_1$ is legal under $FD_2$ as well

**Derivation**: az $FD_1, FD_2, ..., FD_n$ sorozat az $\alpha \to \beta$

(tru)↓ transl. needed levezetése az FD függőségi halmazból, ha

$FD_n = \alpha \to \beta$ (azaz elérhük a célt) és $FD_k \in FD$ (már benne volt) vagy $\exists FD_i, i < k : FD_i \models FD_k$, az Armstrong axiomákkal bizonyíthatóan.

### Armstrong axioms:

- sound: every FD obtained by derivation is consequence of the original one
- complete: if it is consequence, a derivation exists.

### Closure of FD set:

$F^+$ contains all FD's implied by F. (derivable)

### Keys as FD's:

$\beta \subseteq R$ is SUPERKEY, if $\beta \to R$

$\alpha \subseteq R$ is CANDIDATEKEY, if $\alpha$ is super, but $\not\exists \gamma \subset \alpha :$
$\gamma \to R$ is true.

### Closure of attribute set: $\alpha^+ = \{A \mid \exists FD \in F^+ : FD = \alpha \to A\}$

(tru) "azok az attribútumok, amit $\alpha$ meghatároz"

"Algorithm": $\alpha^+(0) = \alpha$; $\alpha^+(i+1) = \alpha^+(i) \cup A$ if $\exists Y \in \alpha^+(i)$
$Y \to A$.

### Theorem: $\alpha \to \beta$ is derivable from F iff $\beta \in \alpha^+$

"$\Rightarrow$" if $\beta \in \alpha^+$ : algorithm constuctively gives derivation

"$\Leftarrow$" if $\alpha \to \beta$ is derivable : by definition $\beta \in \alpha^+$

## Normalization

__Normalization__ : process of
    a.) determining the normal form of a relation
    b.) applying decomposition methods to achieve that.

__Anomalies__ (of poorly designed systems) :
    a.) update anomaly: caused by uncontrolled redundancy
    b.) insert anomaly : not possible to stroe information unless another information is stored as well.
    c) delete anomaly: cannot delete information without deleting another information as well.

__Normal forms__ (shortly) :



> 1NF : all attributes must have atomic domain
>
> 2NF : nonprime attribútum parciálisan nem függhet a
>      CK - től   + 1NF teljesül
>
> 3NF : nonprime attr. nem függhet tranzitív módon
>      CK - től   + 2NF
>
> BCNF : az attribútumok csak SK-től függhetnek

__Total dependency__. in the FD $\alpha \to \beta$ attribute set
   $\beta$ is totally depenten on $\alpha$ if $\nexists \delta \subset \alpha : \delta \to \beta$

__Partial dependency__: in the FD $\alpha \to \beta$ attribute set
   $\beta$ is partially dependent on $\alpha$ if $\exists \delta \subset \alpha : \delta \to \beta$

__Prime attribute__. attribute involved in (at least one) CK.

__Nonprime attribute__ : no CK contains it.

$\boxed{2NF}$ : ~~minim~~ all nonprime attributes **totally depend** on any CK.

※     $\alpha \to \beta$ **violates** the 2NF, if $\beta$ is nonprime, and $\exists \delta \subset \alpha : \delta \to \beta$.

~~3NF~~

Transitive dependency: $\gamma$ depends transitively on a set of attribute $\alpha$ if $\exists \beta : \alpha \to \beta$ and $\beta \to \gamma$.

$\boxed{3NF}$ :

   DEF 1 : A relation is in 3NF, if it is in 2NF and non-prime attributes do not depend ~~partially~~ transitively on any CK.

    $\forall \, \alpha$ nonprime $\not\exists \beta$ CK, $\gamma$ : $\beta \to \gamma$ and $\gamma \to \alpha$.

   DEF 2 : A relation is in 3NF, if for each nontrivial FD in $F^+$ in form of $\alpha \to \beta$ :

     a.) $\alpha$ is superkey    OR :

     b.) $\beta$ is prime attribute

Definition equivalency :

DEF 1 $\to$ DEF 2 :

   if $\alpha \to \beta$ violates DEF 2, ~~two~~ One possible cases.

   a.) $\alpha$ not superkey, $\beta$ nonprime

     if $\alpha$ is not superkey, then :

      i) $\alpha \subseteq CKey \Rightarrow \alpha \to \beta$ is partial dependency of nonprime attribute $\Rightarrow$ not 2NF

      ii) $\alpha \not\subseteq CKey \Rightarrow CKey \to \alpha \to \beta$ is transitive dependency

$-10-$

DEF 2 ⟷ DEF 1 :

if α→β violates DEF 1:
1.) α→β violates 2NF, meaning β is nonprime and
   α is not a CK ⇒ β nonprime, α not superkey
2.) α→β violates transitivity, meaning ~~β~~
   i.) β⊂α ⇒ ~~α~~ β not superkey & γ nonprime
   ii) β⊄α ⇒ β not superkey (??)

---

BCNF: attributes may depend just on superkeys !

[Decomposition methods and conditions]

Decomposition of $R=(A_1, A_2, ..., A_n)$ is a set of
relations $R_1(A_{11}, A_{12}, ..., A_{1i})$
$$R_2(A_{21}, A_{22}, ..., A_{2j})$$
$$\vdots$$
$$R_k(A_{k1}, A_{k2}, ..., A_{kn})$$
such that $\cup A_{ij} = \{A_1, A_2, ..., A_n\}$
and $r_k = \pi_{A_{k1}, A_{k2}, ..., A_{kn}}(r)$.

Lossy decomposition $r \subset r^* = r_1 \bowtie r_2 \subset$ ⇜

Lossless decomposition .

Theorem : if $R_1, R_2, ..., R_k$ is a decomposition of $R$,
   then $r \subseteq r_1 \bowtie r_2 \bowtie ... \bowtie r_k$ is true.
Definition : decomposition is lossless, if $r = r_1 \bowtie ... \bowtie r_k$

<u>Theorem of lossless decomposition</u>: decomposition is
lossless iff at least one holds:

$$R_1 \cap R_2 \rightarrow R_1$$
$$R_1 \cap R_2 \rightarrow R_2$$

the common attribute is
SK in one of the parts

<u>Lossless 2NF decomposition</u>: if $\alpha \rightarrow \beta$ violates the
2NF, decompose to : $R_1 (\alpha, \beta)$
$$R_2 (R-\beta)$$

What ensures its really lossless: $\alpha$ is SK in $R_1$.

<u>Lossless 3NF/BCNF decomposition</u>: always exists

1. Find CK's for R
2. If $\alpha \rightarrow \beta$ in $F^+$ corresponds to 3NF/BCNF
3. Decompose relation. (??)
... details later...

<u>Projection of FD sets</u>:
$$\pi_{R_i}(F) = \{X \rightarrow Y \mid X \rightarrow Y \in F^+ \text{ & } X \cup Y \subseteq R_i \}$$

<u>Decomposition of schemas</u>:

$S = (R; F)$ can be decomposed into
$S_1 = (R_1; F_1) \quad S_2 = (R_2; F_2) \dots \quad S_k = (R_k, F_n)$
Such that:
$$R = \overset{e}{\underset{i=1}{U}} R_i$$
$$F_i = \pi_{R_i}(F)$$

<u>Equivalency of FD sets</u>: E and F are equivalent
iff $E^+ = F^+$, meaning $E \models F$ and $F \models E$.

# Dependency preserving decomposition:

$S \to S_1, S_2, \ldots, S_k$ is dependency preserving, if

$F$ is equivalent to $\overset{k}{\underset{i=1}{\cup}} F_i$

# Ulmans minimal cover (FD set)

1.) single attribute on the right side (decompos. rule)

2) left - reduced : no partial dependency

3.) non-redundant : removing $X \to Y$ ruins equivalency

# Lossless dependency-preserving 3NF decomposition

1.) group dependencies based on the left side

2.) decompose everything based on groups

3.) add relation for the candidate keys

4.) decomposition ok.

What ensures losslessness?

$R_1 \cap R_2 \to R_2$

What ensures dependency preservation: definition of the minimal cover.

# Why does lossless dep.pres BCNF always exist?

$R(City, Street, Zip code)$ $F = \{CS \to Z ; Z \to C\}$

Nonprime : —        Prime : $C, S, Z$.

1NF, 2NF, 3NF ok.

BCNF: $Z$ is not a superkey.

$F^- = \{CS \to Z, Z \to C\}$

$R_1 = (Z, C)$  $Z \to C$ OK

$R_2 = (C, S)$  trivial

# Database recovery systems

__Aim__ : restore database to a known consistent state.

__Properties of database__ :

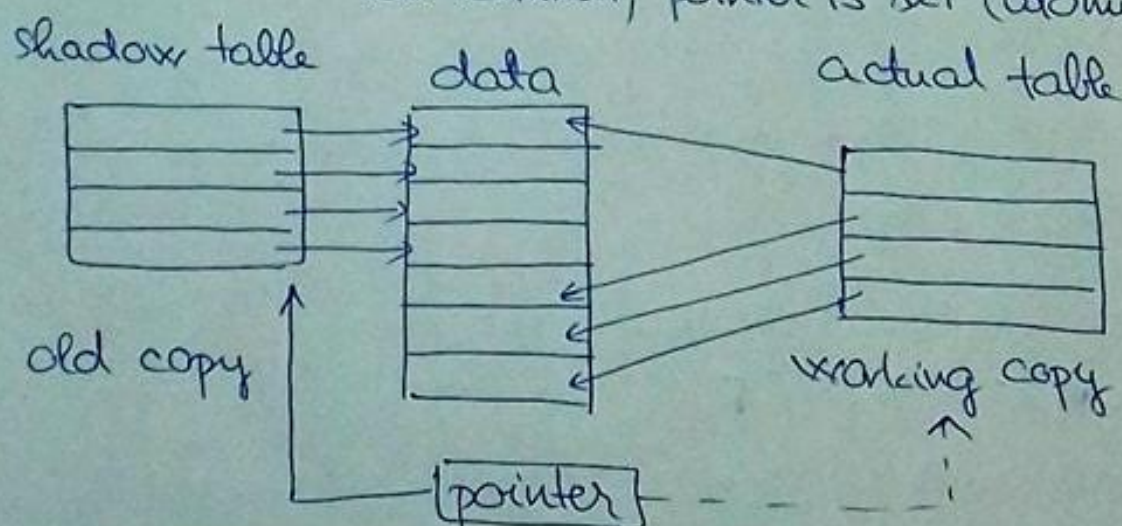Atomicity : transaction completes or does not touch data

Consistency : } later in transaction manager
isolation :

Durability : data is not lost in case of system failure
(like hardware error, power loss, ...)

Types of recovery : LOG based or SHADOW based

__SHADOW based__ : transactions work on a copy

on commit, pointer is set (atomic op).



Shadow table — data — actual table

old copy — working copy

pointer

__LOG based__ :

Deferred modification model (~~UNDO~~ REDO) : uncommitted
transactions are not allowed to touch database

Immediate modification (UNDO)

REDO/UNDO : no idea when storage is written

CHECKPOINT : data of committed transactions is safe

Cascading rollback caused by dirty reading
Solution : immediate modification
deferred + RIGOROUS 2-PHASE PROTOCOL