

```
WITH nameTable1 AS ( SELECT * FROM myTable1 ),
nameTable2 AS ( SELECT * FROM nameTable1 )
SELECT * FROM nameTable2;
```

Rekursiv: zuerst Initialisierung, dann Rekursiver Teil

```
WITH RECURSIVE unter (persnrn, name, chef) AS (
  SELECT A.persnrn, A.name, A.chef FROM angestellter A
  WHERE A.chef = 1010
  UNION ALL
```

```
SELECT A.persnr, A.name, A.chef FROM angestellter A
INNER JOIN unter B ON B.persnr = A.chef
)
SELECT * FROM unter ORDER BY chef, persnr;
```

Views

Sicherheit: Irrelevante Daten für bestimmte Nutzer entfernen.
View ist eine virtuelle Tabelle basierend auf andere Tabellen oder Views. Daten werden zur Ausführzeit aus Tabellendaten hergeleitet. Man kann auch Queries damit vereinfachen.
CREATE VIEW AngPublic (Persnr, Name, Tel, Wohnort) AS
SELECT Persnr, Name, Tel, Wohnort FROM Angestellter;

Updateable Views:

- Wenn nicht enthalten:
 - Join, Set-Operationen
 - Gruppen-Funktionen (min, max)
 - GROUP BY, CONNECT BY, START WITH
 - DISTINCT

Weitere Views

- Materialized Views:
 - Gespeicherte Views
 - nicht automatisch aktualisiert

```
CREATE MATERIALIZED VIEW name... ;
REFRESH MATERIALIZED VIEW name...;
```

- Row-Level Security (RLS)
 - Eine Art "System-Views"
 - Nur User mit entsprechendem Lese- und Schreibrecht («Policy»)

Temporäre Tabellen (CREATE TEMPORARY TABLE):

- Werden gelöscht (dropped) am Ende einer Session oder Transaction
- Andere «permanente» Tabellen mit gleichem Namen sind nicht sichtbar

3-Wertige Logik und Null Werte

Logik: TRUE, FALSE, UNKNOWN

NULL-Werte werden bei Aggregatfunktionen nicht mit verwendet

More operations with NULL which might be surprising:

2 + NULL	=> NULL	NULL = NULL	=> NULL
NULL + 10	=> NULL	2 = NULL	=> NULL
CONCAT('hi', NULL)	=> NULL	2 != NULL	=> NULL

Datenbank-Sicherheit (DCL – Data Control Language)

- System-Sicherheit (Allgemeiner Datenschutz)
 - Authentisierung: Benutzer muss sich identifizieren mit Passwort und Benutzername
 - Benutzerprivilegien (Autorisierung): Darf Nutzer Das
 - Kontrolle von System-Ressourcen, Auditing, Transportsicherheit
- Daten-Sicherheit (Datenschutz DBMS)
 - Zugriffskontrolle zur Verhinderung von nicht autorisiertem Zugriff auf Datenobjekte
 - Auditing von Zugriffsoperationen
 - DCL ist nicht standardisiert
- Massnahmen gegen Exploids, SQL-INJECTION:
- User-Input eingrenzen (Typ prüfen (z.B numeric), Escapen von nicht numerischen Input
- Datenzugriff abstrahieren (Stored Procedures)
- Spezieller Nutzer mit limitierten Rechte (Superuser)
- Benutzerverwaltung:
 - Jedem User sind Rechte zugeordnet für die Datenbankoperationen (Login/Connect)
 - Jedem User sind Rechte zugeordnet für die Verwaltung von Datenbankobjekte
- Schema:
 - Schema fassen Datenbank-Objekte zusammen in einer bestimmten Datenbank
 - Datenbank kann n-Schemas haben
 - Default: public
- Rolle (ROLE):
 - Oberbegriff für User oder Gruppen
 - Rollen gelten für den ganzen Cluster (über alle Datenbanken)
 - Rolle kann n-Schemas besitzen
- Benutzer: Rolle mit LOGIN
 - CREATE USER ... ist identisch mit CREATE ROLE ... WITH LOGIN

Gruppe: ROLE ohne LOGIN

Privilegien:

- Systemprivilegien
 - Erlauben den Zugriff auf die Datenbank-Operationen
 - CREATEDB, CREATEROLE
- CREATE ROLE user CREATEDB NOCREATOROLE;
ALTER ROLE user WITH CREATEROLE;
- Datenprivilegien
 - Erlauben den Zugriff auf die Datenobjekte
 - Grantor gewährt Privilegien
 - Grantee erhält Privilegien
 - GRANT object_privileg ON object TO (user|group|PUBLIC);
 - object_priv: - SELECT, INSERT, UPDATE, DELETE, REFERENCES, TRIGGER, ALL
 - REVOKE – Gegenteil von GRANT (Entfernt die vergebenen Privilegien)
 - Gruppe - User zu Gruppe zuordnen: GRANT gruppe TO user;
 - Gruppe sind globale Objekte, nicht in einem Schema!
 - Nur Objektprivilegien

Transaktionen

Pro Session maximum 1 Transaktion
NESTED Transaktion nicht unterstützt.

Nutzen von Transaktionen:

Fault Tolerance: Bei Servercrash kann operation wiederholt werden oder wird ganz gecancelt

Concurrency: Isolation der Transaktionen, Parallelität wird ermöglicht

ACID

- Atomicity:** Vollständig oder gar nicht
- Consistency:** Konsistenter Zustand bleibt erhalten
- Isolation:** Transaktion soll von anderen Isoliert sein
- Durability:** Alle Änderungen sind persistent

```
BEGIN TRANSACTION; -- Kurznotation -> BEGINN;
COMMIT TRANSACTION; -- Kurznotation -> COMMIT;
ROLLBACK TRANSACTION; -- Kurznotation -> ROLLBACK;
SAVEPOINT blabla;
ROLLBACK TO blabla;
RELEASE blabla;
COMMIT;
```

- Commit Resultate
 - Success
 - Änderungen atomar und durable gespeichert
 - Failure
 - Alle temporären Änderungen werden abgebrochen (abort)
- Gründe für Abort
 - Explizit durch ROLLBACK oder ABORT
 - Unzulässige Verzahnung mit anderen nebenläufigen Transaktionen, Deadlock
 - Applikationsabbruch, Systemabsturz, Fehler

Serialisierbarkeit

wenn, nebenläufige Ausführung, gleich wie serielle Ausführung Muss azyklisch sein -> Serialisierbarkeit keine Schlaufen

Beispiel:

- T1 = r1(x)w1(x)r1(y)c1 und T2 = r2(x)w2(x)r2(y)w2(y)c2
- S = r1(x)r2(x)w1(x)w2(x)r2(y)r1(y)w1(y)w2(y)c2c1

- Konfliktpaare :
 - r1(x)→w2(x), r2(x)→w1(x), w1(x)→w2(x)
 - Von w(x) zu allen ..(x) zeichnen, Pfeile links nach rechts
 - Immer mit einem Write vorhanden

Serialisierungsgraphen

- Jedes Konfliktpaar eine Verbindung (Ausnahme T1 zu T1)
- T1 → T2 Zyklus → Nicht serialisierbar
- Topologische Sortierung (Halbordnung) bestimmt Commit-Reihenfolge

Implementation der Isolation

- Pessimistische Verfahren
 - Sperrprotokolle
 - Besser bei hoher Konflikt-Wahrscheinlichkeit
- Optimistische Verfahren
 - Konfliktbehebung im Nachhinein
 - Besser bei kleiner Konfliktwahrscheinlichkeit

Locking

Locks auf verschiedene Granulitäten (Table, Table Range, Row, Item). Desto kleiner, desto grössere Parallelität. Psql meist ROW.

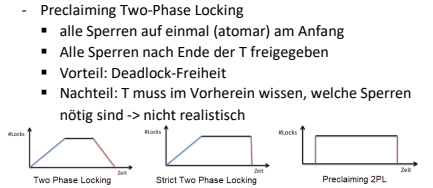
- Exclusive Lock (X)

- Für Schreibe- oder Lesezugriffe
- Nur eine Transaktion *lock*
- Shared Lock (S)
 - Nur für Lese-Zugriffe
 - Mehrere Transaktionen *slock*

Locking garantiert keine Serialisierbarkeit → zu frühes Unlock.

- Two Phase Locking (2PL)
 - Phase 1(Growing Phase): Objekt sperre vor Zugriff
 - Phase 2(Shrinking Phase): nach einem unock, kein lock
- Strict Two Phase Locking
 - Alle Sperren nach Ende der T freigeben
- Vorteile zum 2PL
 - kein Cascading Rollback
 - kein unklarer Beginn der Shrinking Phase
- Nachteile
 - Deadlocks sind möglich
 - beschränkt Parallelität unnötig ein

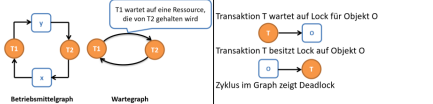
- Preclaiming Two-Phase Locking
 - alle Sperren auf einmal (atomar) am Anfang
 - Alle Sperren nach Ende der T freigeben
 - Vorteil: Deadlock-Freiheit
 - Nachteil: T muss im Vorheren wissen, welche Sperren nötig sind -> nicht realistisch



Deadlock Szenario

Gegenseitige Locks. Wird entweder durch Timeout (poor man solution) oder durch die Erkennung (Abbruch einzelner Transaktionen abgebrochen).

Für Analyse dient Betriebsmittelgraph und Wartegraph:



- Deadlock-Auflösung
 - Scheduler erkennt Zyklus im Wartegraph
 - Opfersuche bei Zyklus
 - Wähle eine Transaktion im Zyklus und breche diese ab
 - Minimierung des Rücksetzaufwands, T mit Wenigsten sperren
 - Maximierung frei gewordene Ressourcen: T mit meisten Sperren
 - Vermeidung der Starvation: Nicht immer dieselbe T rollbacks
 - Mehrfache Zyklen zerstören
 - Cascading Rollbacks: Implizite Rollbacks von weiteren T

Isolation

SERIALISED am besten, doch Parallelität limitiert, deshalb wird Effizienz mit schwächeren Levels gesteigert, auf kosten der Korrektheit. Fehler schwer nachvollziehbar.

4 Level nach ANSI SQL-92 Standard

- READ UNCOMMITTED
 - Lesezugriffe nicht synchronisiert (keine read-lock)
 - Read ignoriert jegliche Sperren
- READ COMMITTED
 - Lesezugriffe nur kurz/temporär synchronisiert (default)
 - setzt für gesamte T Write-Lock, Read-Lock nur kurzfristig
- REPEATABLE READ
 - einzelne Zugriffe ROWS sind synchronisiert
 - Read und Write Lock für die gesamte T
- SERIALIZABLE
 - Vollständig (korrekte) Isolation ACID

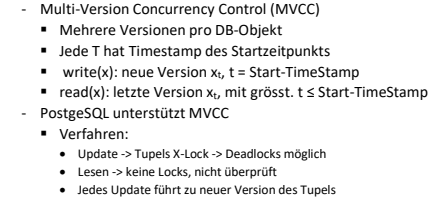
	Read Uncommitted	Read Committed	Repeatable	Serializable
Dirty Write	möglich	möglich	möglich	unmöglich
Dirty Read	möglich	unmöglich	unmöglich	unmöglich
Lost Update	möglich	möglich	unmöglich	unmöglich
Fuzzy Read	möglich	möglich	unmöglich	unmöglich
Phantom	möglich	möglich	möglich	unmöglich
Read Skew	möglich	möglich	unmöglich	unmöglich
Write Skew	möglich	möglich	möglich	möglich
Deadlock			möglich	unmöglich
Cascading Rollback				unmöglich

- Dirty Read – Daten Lesen von anderer nicht committed T
- Fuzzy Read – Lese gleiche Daten mehrmals, sehe aber andere Werte, gelesene Daten ändern sich durch andere T

- Phantom – gleiche Selects entdecken neue/gelöschte Rows
- Serializable: Können Leser, Schreiber blockieren
- Read Committed: Können Schreiber, Schreiber blockieren
- Es gilt immer höchste Isolation
 - SET TRANSACTION ISOLATION LEVEL ...;

Alternative: Optimistische verfahren ohne Locks.

- Jede Transaktion sieht Snapshot zu Start-Zeitpunkt
- Bei Änderung im Commit prüfen, dass Objekte unverändert sind, wie zum Snapshot-Zeitpunkt, sonst Rollback



- Multi-Version Concurrency Control (MVCC)
 - Mehrere Versionen pro DB-Objekt
 - Jede T hat Timestamp des Startzeitpunkts
 - write(x): neue Version x_n, t = Start-TimeStamp
 - read(x): letzte Version x_n, mit grösst. t ≤ Start-TimeStamp
- PostgreSQL unterstützt MVCC
- Verfahren:
 - Update -> Tupels X-Lock -> Deadlocks möglich
 - Lesen -> keine Locks, nicht überprüft
 - Jedes Update führt zu neuer Version des Tupels

	Garantiert Serialisierbar	Keine Deadlocks	Keine Cascading Rollbacks	Keine Konflikt-Rollbacks	Hohe Parallelität	Realistisch (ohne Voranalyse)
Two-Phase Locking	✓	✗	✗	✓	✗	✗
Strict 2PL	✓	✗	✗	✓	✗	✓
Preclaiming 2PL	✓	✓	✓	✓	✓	✗
Validation-Based	✓	✓	✗	✗	✓	✓
Timestamp-Based	✓	✓	✗	✗	✓	✓
Snapshot Isolation	✗	✗	✗	✗	✓	✓
SSI	✓	✗	✓	✗	✓	✓

! Deadlock in PS

Log Files

- Write-Ahead Log (WAL)
 - Änderung der T in Log Schreiben (Flush)
 - Ebenso Commit in Log schreiben atomar (Flush)
 - Schlussendlich In-Place Updates in DB
- Aufbau: [LSN, TaId, PageId, Redo, UNDO, PrevLSN]
 - LSN (Log Sequence Number): Eindeutig monoton ansteigend
 - TransaktionsID: T-Kennung, die die Änderung durchgeführt hat
 - Page: Page-Kennung, auf der die Änderung vollzogen wurde
 - Redo: Absoluter Wert nach der Änderung
 - Undo: Absoluter Wert vor Änderung
 - PrevLSN: vorhergehende Log-Eintrag der jeweilige T (effizient)
- Nach Absturz Recovery: Undo nach Log für alle COMMIT T
 - Winners: Committed T müssen wieder REDO werden
 - Looses: Nicht abgeschlossene T werden UNDO werden

Backup

- Planung
 - Maximale Ausfallzeit, Zeit für Recovery
 - Was soll gesichert werden? (Read-Only via Import=)
 - Wann soll gesichert werden?
 - Voller Backup (Log- und Datenfiles)
 - Inkrementelles Backup (Logfiles)
 - Export (Daten)
 - Online/Offline Backup
- Speicherstrukturen
 - Spiegelung des Log-Files, durch DBMS, Durch OS/Hardware
- Backup-arten
 - Logischer Backup – SQL DUMP (pg_dump)
 - Blockiert keine Schreibende oder Lesende T
 - Für Mitteltgrosse Datenmengen
 - Interkompatibel mit neuen PG-Versionen, andere Maschinen
 - Physisches Backup – File System
 - Datenbank muss gestoppt werden

- Schneller als Logische Backup
- Passst nur zu deselben MAJOR Version
- Andere: Cloud, ContinuosArchiving, Snapshot, Agent

Indexe

- Data Pages – Heap
- Index Pages – Suchbaum (B-Baum)
 - Wurzel-Knoten(Root-Node), Knoten(Node), Blätter(Leaf)
 - Primär-Index immer unique

Arten:

- ISAM (Index-Sequential Access Method) [Historisch zuerst]
 - Einfügen und Suchen: einfach und schnell
 - Aktualisieren: schlecht
 - Daten werden über die Indexspalte aufsteigend sortiert
- B-Bäume
 - B für Balancierter Baum -> Überall gleich viele Stufen
 - Gecusterter Index: Nachbar Leafs referenzieren auch auf Nachbarden
 - Grad k > Leaf & Node[k; 2k] Einträge, Root [1; 2k]
 - Garantierte Mindestauslastung beträgt 50% (Real:49.99% ->Root evt 1 Eintrag)
- Einfügen:
 - Suche nach dem Schlüssel -> Suche endet in den Blättern
 - Füge den Schlüssel dort ein
 - Ist der Knoten(Blatt) überfüllt?
 - 3.1 Neuer Knoten mit dem mittleren und der rechts liegenden Einträge
 - 3.2 Kleinste Eintrag im neuen Knoten in den Vaterknoten
 - 4. Ist der Vaterknoten (jetzt) überfüllt? -> Analog zu den Normalen Knoten -> Root überfüllt? -> Baum wächst
- B+-Bäume
 - Ein B-Baum: Daten oder die Referenz nur in den Blättern (unterste Stufe)
 - Daten in den Blättern sehr selten
- Hash
 - Ordnet Key zu Einträgen
 - Problem: Overflow
 - Spezielle Speicherstrukturen bei PostgreSQL

Indexe in PostgreSQL:

- B-Tree (Default bei PostgreSQL)
 - Universell: Bereichsabfragen, Vergleiche, Mustersuche
- Hash (Ab PG Version 9.3 nicht nutzen)
 - B-Tree in fast allen Aspekten besser
- GIST (Generalized Search Tree)
 - Baumartige Struktur
 - Range/Containment Search, K-Nearest-Neighbour Search
 - Für geometrische Datentypen, Volltextsuche
- GIN (Generalized Inverted Index)
 - Geeignet für Arrays
 - speichert effizient Duplikate
 - VS GIST: erstellen, update langsamer, Zugriff schneller
- BRIN (Block Range Index)
 - Speichert min/max-Werte als «Blöcke»
 - Gut für «Range Search»
 - natürlich benachbarte oder sortierte Daten
 - Nutzt natürliche Nähe, z.B. Einfügedatum, PLZ
 - kleiner Disk-Verbrauch
- Index-Variationen
 - Zusammengesetzter Index
 - Index über mehrere Attribute/Kolonnen
 - CREATE INDEX table_idx ON table (col1,col2);
 - Index mit INCLUDE
 - CREATE INDEX table_idx ON table (col1) INCLUDE (col2);
 - Partieller Index: Queries beziehen sich auf eine Selektion
 - CREATE INDEX t_idx ON table (c1) WHERE c2 IS NOT NULL;
 - Funktionaler Index
 - Index mit Funktion / Ausdruck
 - CREATE INDEX table_idx ON table (function (col1));

Optimierungen (=> EXPLAIN ANALYZE)
Ist Query I.O.? / ist Schema i.O.? / DB-Konfiguration HW angepasst? / DB-Statistik nachgeführt? / richtige Indexe vorhanden? / Index geclustert? / DB-Architektur adäquat?
Logisch: Abfrage umformulieren, Physisch: Algorithmus ändern