

11. Tabellenoperationen – Implementierung

- Ziele¹
- Systematische Entwicklung von relationalen Verarbeitungskonzepten für eine oder mehrere Tabellen
- Realisierung von Planoperatoren

• Operationen der Relationenalgebra

- unäre Operationen: π, σ
- binäre Operationen: $\bowtie, \times, \div, \cap, \cup, -$

➔ SQL-Anfragen enthalten logische Ausdrücke, die auf die Operationen der Relationenalgebra zurückgeführt werden können. Sie werden in Zugriffspläne umgesetzt. Sog. Planoperatoren implementieren diese logischen Operationen

• Planoperatoren auf einer Tabelle

• Selektion

• Operatoren auf mehreren Tabellen

• Verbundalgorithmen

- Nested-Loop-Verbund
- Sort-Merge-Verbund
- Hash-Verbund
(classic hashing, simple hash join, hybrid hash join)
- Nutzung von typübergreifenden Zugriffspfaden
- verteilte Verbundalgorithmen

• weitere binäre Operationen (Mengenoperationen)

¹ Mitschang, B.: Anfrageverarbeitung in Datenbanksystemen – Entwurfs- und Implementierungskonzepte, Reihe Datenbanksysteme, Vieweg-Verlag, 1995

Planoperatoren auf einer Tabelle

• Selektion – Allgemeine Auswertungsmöglichkeiten:

- direkter Zugriff über ein gegebenes TID, über ein Hash-Verfahren oder eine ein- bzw. mehrdimensionale Indexstruktur
- sequentielle Suche in einer Tabelle
- Suche über eine Indexstruktur (Indextabelle, Bitliste)
- Auswahl mit Hilfe mehrerer Verweislisten, wobei mehr als eine Indexstruktur ausgenutzt werden kann
- Suche über eine mehrdimensionale Indexstruktur

• Projektion

wird typischerweise in Kombination mit Sortierung, Selektion oder Verbund durchgeführt

• Modifikation

- Änderungen sind in SQL mengenorientiert, aber auf eine Tabelle beschränkt
- INSERT, DELETE und UPDATE werden direkt auf die entsprechenden Operationen der Speicherungsstrukturen abgebildet
- „automatische“ Abwicklung von Wartungsoperationen zur Aktualisierung von Zugriffspfaden, zur Gewährleistung von Cluster-Bildung und Reorganisation usw.
- Durchführung von Logging- und Recovery-Maßnahmen usw.

Planoperatoren für die Selektion

- **Nutzung des Scan-Operators**

- Definition von Start- und Stopp-Bedingung
- Definition von einfachen Suchargumenten

- **Planoperatoren**

1. *Tabellen-Scan (Relationen-Scan)*

- immer möglich
- SCAN-Operator implementiert die Selektionsoperation

2. *Index-Scan*

- Auswahl des kostengünstigsten Index
- Spezifikation des Suchbereichs (Start-, Stopp-Bedingung)

3. *k-d-Scan*

- Auswertung mehrdimensionaler Suchkriterien
- Nutzung verschiedener Auswertungsrichtungen durch Navigation

4. *TID-Algorithmus*

- Auswertung aller „brauchbaren“ Indexstrukturen
- Auffinden von variabel langen TID-Listen
- Boolesche Verknüpfung der einzelnen Listen
- Zugriff zu den Tupeln entsprechend der Trefferliste

- **Weitere Planoperatoren in Kombination mit der Selektion**

- Sortierung
- Gruppenbildung
(siehe Sortieroperator)
- spezielle Operatoren z. B. in Data-Warehouse-Anwendungen zur Gruppen- und Aggregatbildung (CUBE-Operator)

Operatoren über mehrere Tabellen

- **SQL erlaubt komplexe Anfragen über k Tabellen**

- **Ein-Variablen-Ausdrücke:**
beschreiben Bedingungen für die Auswahl von Elementen aus einer Tabelle.
- **Zwei-Variablen-Ausdrücke:**
beschreiben Bedingungen für die Kombination von Elementen aus zwei Tabellen.
- k-Variablen-Ausdrücke werden typischerweise in Ein- und Zwei-Variablen-Ausdrücke zerlegt und durch entsprechende Planoperatoren ausgewertet

- **Planoperatoren über mehrere Tabellen**

Allgemeine Auswertungsmöglichkeiten:

- **Schleifeniteration** (*nested iteration*)

für jedes Element der äußeren Tabelle R_a Durchlauf der inneren Tabelle R_i

- $O(N_a \cdot N_i + N_a)$
- wichtigste Anwendung: *nested loops join*

- **Mischmethode** (merge method)

sequentieller, schritthaltender Durchlauf beider Tabellen R_1, R_2

- $O(N_1 + N_2)$
- ggf. zusätzliche Sortierkosten
- wichtigste Anwendung: *merging join*

- **Hash-Methode** (*hashing*)

Partitionierung der inneren Tabelle R_i . Laden der p Partitionen in eine Hash-Tabelle HT im HSP. „Probing“ der äußeren Tabelle R_a oder ihrer entsprechenden Partitionen mit HT:

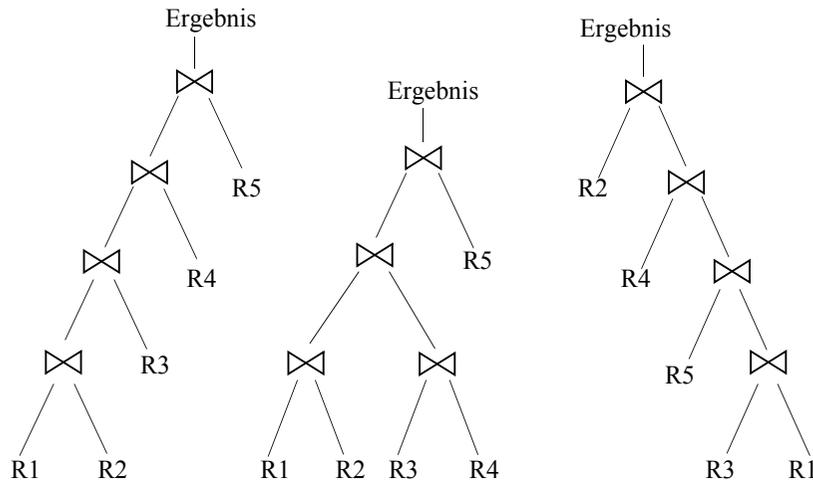
$$\rightarrow O(p \cdot N_a + N_i)$$

Operatoren über mehrere Tabellen (2)

• n-Wege-Verbunde

- Zerlegung in n-1 Zwei-Wege-Verbunde²
- Anzahl der Verbundreihenfolgen ist abhängig von den gewählten Verbundattributen
- maximal **n!** verschiedene Reihenfolgen möglich
- Einsatz von Pipelining-Techniken
- Optimale Auswertungsreihenfolge abhängig von
 - Planoperatoren
 - „passende“ Sortierordnungen für Verbundattribute
 - Größe der Operanden usw.

• Einige Verbundreihenfolgen mit Zwei-Wege-Verbunden (n=5)



• Analoge Vorgehensweise bei Mengenoperationen

² Praxistest (Guy Lohman test for join techniques): Does a new technique apply to joining three inputs without interrupting data flow between the join operators?

Planoperatoren für den Verbund

• Verbund

- satztypübergreifende Operation: gewöhnlich sehr teuer
- häufige Nutzung: wichtiger Optimierungskandidat
- typische Anwendung: Gleichverbund
- allgemeiner Θ -Verbund selten

- **Implementierung der Verbundoperation** kann gleichzeitig Selektionen auf den beteiligten Tabellen R und S ausführen

```
SELECT *
FROM R, S
WHERE R.VA  $\Theta$  S.VA
      AND PR
      AND PS
```

- VA: Verbundattribute
- P_R und P_S: Prädikate definiert auf Selektionsattributen (SA) von R und S

• Mögliche Zugriffspfade

- Scans über R und S (immer)
- Scans über I_R(VA), I_S(VA) (wenn vorhanden)
 - ↳ liefern Sortierreihenfolge nach VA
- Scans über I_R(SA), I_S(SA) (wenn vorhanden)
 - ↳ ggf. schnelle Selektion für P_R und P_S
- Scans über andere Indexstrukturen (wenn vorhanden)
 - ↳ ggf. schnelleres Auffinden aller Sätze

Nested-Loop-Verbund

- **Annahmen:**

- Sätze in R und S sind nicht nach den Verbundattributen geordnet
- es sind keine Indexstrukturen $I_R(VA)$ und $I_S(VA)$ vorhanden

- **Algorithmus für Θ -Verbund:**

Scan über S,
für jeden Satz s, falls P_S :
 Scan über R,
 für jeden Satz r, falls $P_R \text{ AND } (r.VA \Theta s.VA)$:
 führe Verbund aus,
 d. h., übernehme kombinierten Satz (r, s) in die Ergebnismenge.

- **Komplexität: $O(N^*M)$**

- **Nested-Loop-Verbund mit Indexzugriff**

Scan über S,
für jeden Satz s, falls P_S :
 ermittle mittels Zugriff auf $I_R(VA)$ alle TIDs für Sätze mit $r.VA = s.VA$,
 für jedes TID:
 hole Satz r, falls P_R :
 übernehme kombinierten Satz (r, s) in die Ergebnismenge.

- **Nested-Block-Verbund**

Scan über S,
für jede Seite (bzw. Menge aufeinanderfolgender Seiten) von S:
 Scan über R,
 für jede Seite (bzw. Menge aufeinanderfolgender Seiten) von R:
 für jeden Satz s der S-Seite, falls P_S :
 für jeden Satz r der R-Seite,
 falls $P_R \text{ AND } (r.VA \Theta s.VA)$:
 übernehme kombinierten Satz (r, s) in die Ergebnismenge.

Sort-Merge-Verbund

- **Algorithmus besteht aus 2 Phasen:**

- **Phase 1:** Sortierung von R und S nach R(VA) und S(VA)
(falls nicht bereits vorhanden), dabei frühzeitige
Eliminierung nicht benötigter Sätze ($\rightarrow P_R, P_S$)
- **Phase 2:** schritthaltende Scans über sortierte R- und S-Sätze
mit Durchführung des Verbundes bei $r.VA = s.VA$

- **Komplexität: $O(N \log N)$**

- **Spezialfall**

Falls $I_R(VA)$ und $I_S(VA)$ oder verallgemeinerte Zugriffspfadstruktur über R(VA) und S(VA) (Join-Index) vorhanden:

\rightarrow **Ausnutzung von Indexstrukturen auf Verbundattributen:**

Schritthaltende Scans über $I_R(VA)$ und $I_S(VA)$:
für jeweils zwei Schlüssel aus $I_R(VA)$ und $I_S(VA)$, falls $r.VA = s.VA$:
hole mit den zugehörigen TIDs die Tupel,
falls P_R und P_S :
übernehme kombinierten Satz (r, s) in die Ergebnismenge.

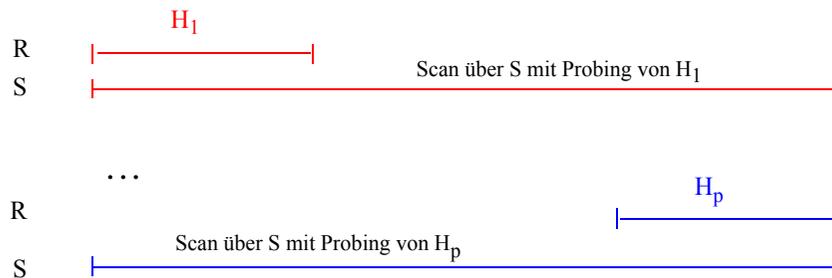
Hash-Verbund

- **Einfachster Fall (classic hashing):**

- *Schritt 1:* Abschnittsweises Lesen der (kleineren) Tabelle R und Aufbau einer Hash-Tabelle mit $h_A(r(VA))$ nach Werten von R(VA) entsprechend den Abschnitten R_i ($1 \leq i \leq p$), so daß jeder der p Abschnitte in den verfügbaren Hauptspeicher paßt und jeder Satz P_R erfüllt
- *Schritt 2:* Überprüfung (Probing) für jeden Satz von S mit P_S ; im Erfolgsfall Durchführung des Verbundes
- *Schritt 3:* Wiederhole Schritt 1 und 2 solange, bis R erschöpft ist.

- **Aufbau der Hash-Tabelle und Probing**

Es erfolgt ein Scan über R; dabei wird die Hash-Tabelle H_i ($1 \leq i \leq p$) der Reihe nach im HSP aufgebaut



- **Komplexität: $O(p \cdot N)$**

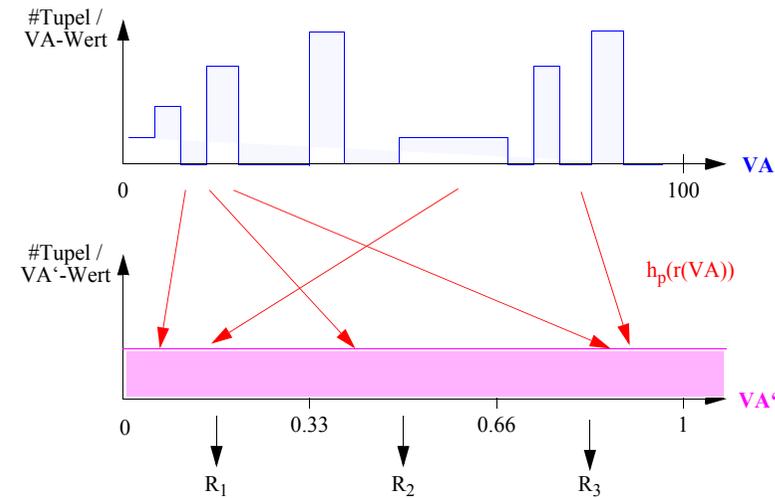
- **Spezialfall:**

R paßt in den Hauptspeicher: eine Partition ($p = 1$)

→ ein Scan über S genügt

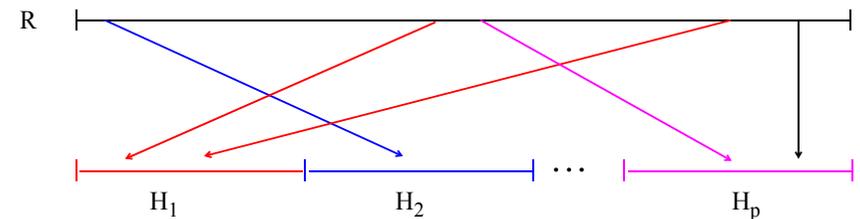
Hash-Verbund (2)

- **Partitionieren von R mit $h_p(r(VA))$**



- **Partitionierung**

- Partitionierung von R in Teilmengen R_1, R_2, \dots, R_p :
Ein Satz r von R ist in R_p , wenn $h(r)$ in H_p ist.



→ Warum ist diese Partitionierung eine kritische Operation?

Welche Hilfsoperationen können erforderlich sein?

Ist für die Partitionierung der Einsatz einer Hash-Funktion notwendig?

- Tabelle S wird mit **derselben Funktion h_p** unter Auswertung von P_S partitioniert

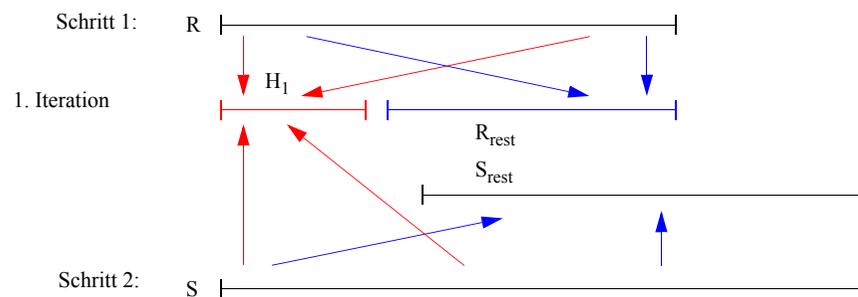
Hash-Verbund (3)

- **Varianten des Hash-Verbundes:**

Sie unterscheiden sich vor allem durch die Art der Partitionsbildung

- **Partitionierungstechnik beim einfachen Hash-Verbund:**

gezeigt am Aufbau und Probing von H_1



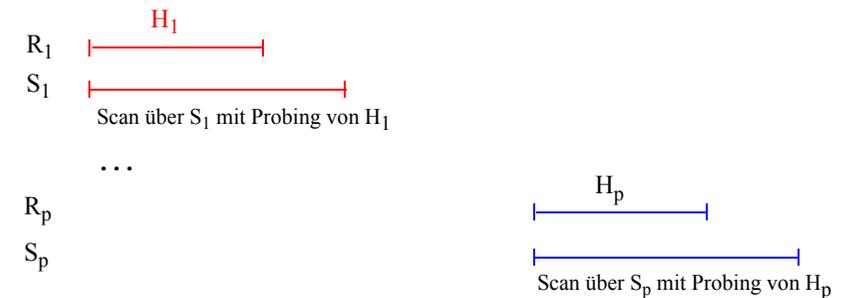
- **Einfacher Hash-Verbund (simple hash join)**

- **Schritt 1:** Führe Scan auf R (kleinere Tabelle) aus, überprüfe P_R und wende auf jedes qualifizierte Tupel r die Hash-Funktion h_p an. Fällt $h_p(r(VA))$ in den gewählten Bereich, trage es in H_1 ein. Andernfalls schreibe r in einen Puffer für die Ausgabe in eine Datei R_{rest} für „übergangene“ r -Tupel.
- **Schritt 2:** Führe Scan auf S aus, überprüfe P_S und wende auf jedes qualifizierte Tupel s die Hash-Funktion h_p an. Fällt $h_p(s(VA))$ in den gewählten Bereich, suche in H_1 einen Verbundpartner (Probing). Falls erfolgreich, bilde ein Verbundtupel und ordne es dem Ergebnis zu. Andernfalls schreibe s in einen Puffer für die Ausgabe in eine Datei S_{rest} für „übergangene“ s -Tupel.
- **Schritt 3:** Wiederhole Schritt 1 und 2 mit den bisher übergangenen Tupeln auf H_1 solange, bis R_{rest} erschöpft ist. Dabei ist die Überprüfung von P_R und P_S nicht mehr erforderlich.

Hash-Verbund (4)

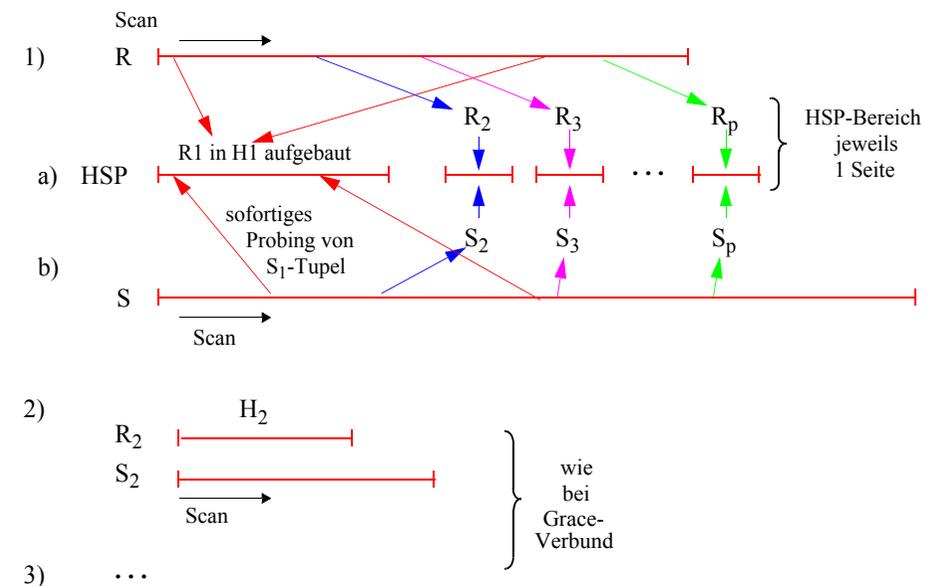
- **Grace-Verbund (grace join)**

- Partitionsbildung findet vor dem Verbund statt
- Partitionen R_i und S_i sind in Dateien zwischengespeichert
- Aufbau von H_i im HSP mit R_i und Probing mit S_i



- **Hybrider Hash-Verbund (hybrid hash join)**

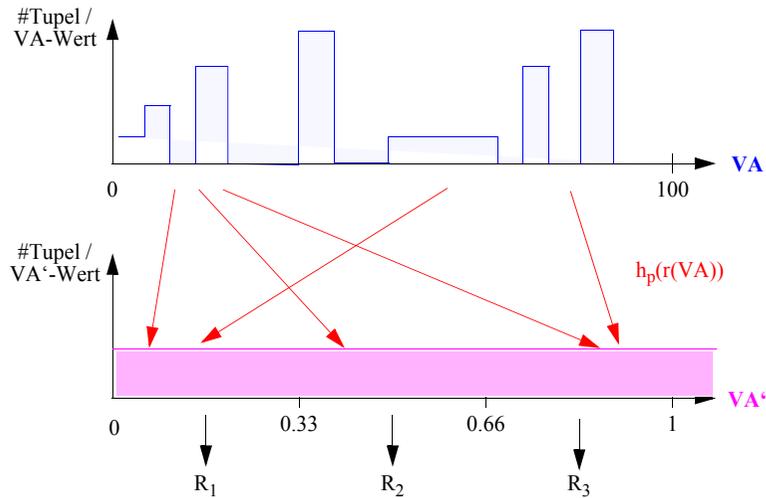
- optimiert das Verfahren dadurch, daß parallel zur Partitionsbildung Aufbau und Probing von H_1 erfolgt



Hash-Verbund – Beispiel

I. Partitionieren

a) Partitionieren von R mit $h_p(r(VA))$



b) Partitionieren von S mit $h_p(s(VA))$

II. Verbund

- 1) R_1 $\xrightarrow{H_1}$ in M mit $h_H(r(VA))$
 $VA': 0.0 - 0.33$
 S_1 $\xrightarrow{VA': 0.0 - 0.33}$
 Lesen, Probing mit $h_H(s(VA))$
- 2) R_2 $\xrightarrow{H_2}$
 $VA': 0.34 - 0.66$
 S_2 $\xrightarrow{VA': 0.34 - 0.66}$
- 3) R_3 $\xrightarrow{H_3}$
 $VA': 0.67 - 1.0$
 S_3 $\xrightarrow{VA': 0.67 - 1.0}$

Nutzung von typübergreifenden Zugriffspfaden

• Verbund über Link-Strukturen

Ausnutzung hierarchischer Zugriffspfade für den Gleichverbund

Scan über R (Owner-Tabelle),

für jeden Satz r, falls P_R :

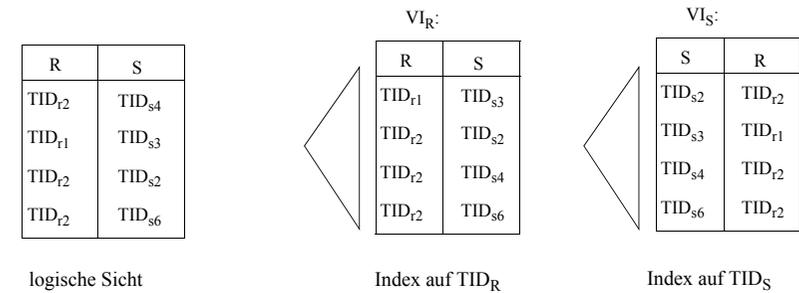
Scan über zugehörige Link-Struktur $L_{R-S}(VA)$,

für jeden Satz s, falls P_S :

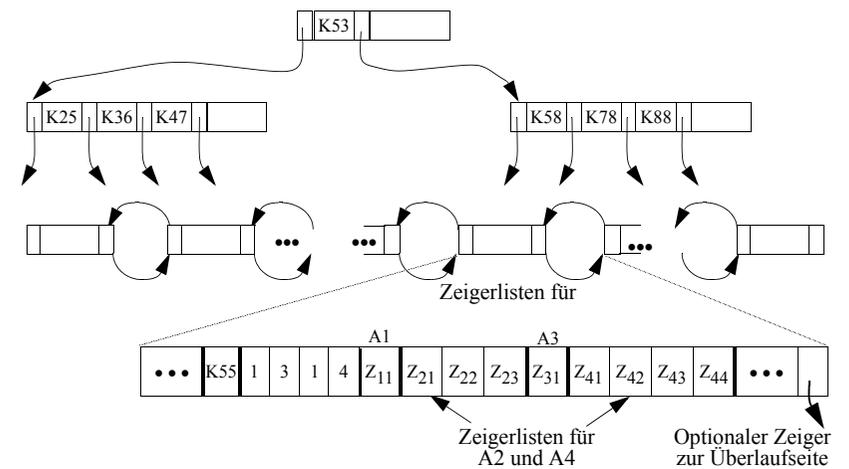
übernehme kombinierten Satz (r, s) in die Ergebnismenge.

• Weitere Verfahren

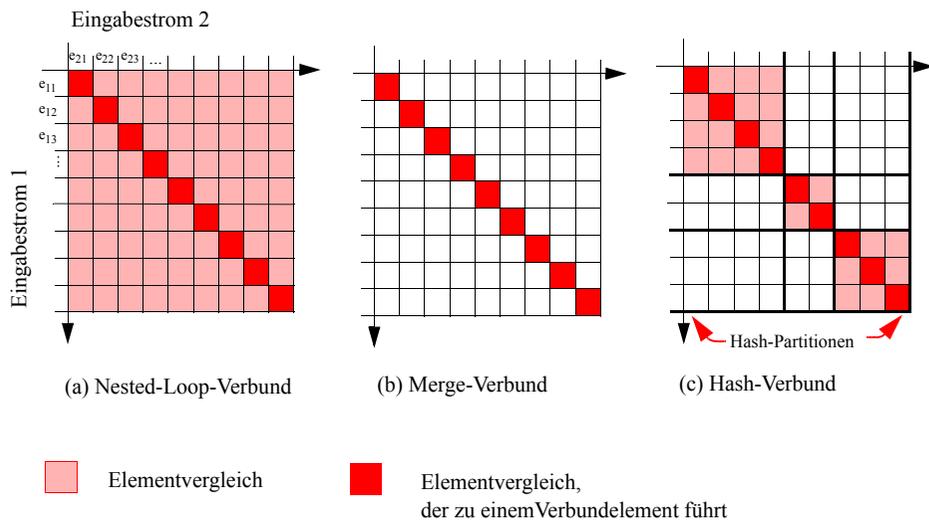
- **Verbundindexe**, die für bestimmte Θ -Verbunde eingerichtet sind



- Nutzung von **verallgemeinerten Zugriffspfadstrukturen**



Verbundalgorithmen – Vergleich



- **Nested-Loop-Verbund** ist immer anwendbar, jedoch ist dabei stets das vollständige Durchsuchen des gesamten Suchraums in Kauf zu nehmen.
- **Merge-Verbund** benötigt die geringsten Suchkosten, verlangt aber, daß die Eingabeströme bereits sortiert sind. Indexstrukturen auf beiden Verbundattributen erfüllen diese Voraussetzung. Sonst reduziert das Sortieren beider Tabellen nach den Verbundattributen den Kostenvorteil in erheblichem Maße. Ein Sort-Merge-Verbund kann dennoch zusätzliche Vorteile besitzen, falls das Ergebnis in sortierter Folge verlangt wird und das Sortieren des großen Ergebnisses aufwendiger ist als das Sortieren zweier kleiner Ergebnismengen.
- Beim **Hash-Verbund** wird der Suchraum partitioniert. In Bild c ist unterstellt, daß die gleiche Hash-Funktion h auf die Tabellen R und S angewendet worden ist. Die Partitionsgröße (bei der kleineren) Tabelle richtet sich nach der verfügbaren Puffergröße im Hauptspeicher. Eine Verkleinerung der Partitionsgröße, um den Fall b anzunähern, verursacht höhere Vorbereitungskosten und ist deshalb nicht zu empfehlen.

Verbundalgorithmen in verteilten DBS

• Problemstellung:

- Anfrage in Knoten K, die einen Verbund zwischen (Teil-) Tabellen R am Knoten K_R und (Teil-) Tabelle S am Knoten K_S erfordert
- Festlegung des Ausführungsknotens: K, K_R oder K_S

• Bestimmung der Auswertestrategie

- Sende beteiligte Tabellen vollständig an einen Knoten und führe lokale Verbundberechnung durch („*ship whole*“““
 - minimale Nachrichtenanzahl
 - sehr hohes Übertragungsvolumen
- Fordere für jeden Verbundwert der ersten Tabelle zugehörige Tupel der zweiten an („*fetch as needed*“““
 - hohe Nachrichtenanzahl
 - nur relevante Tupel werden berücksichtigt
- Kompromißlösung:
Semi-Verbund bzw. Erweiterungen wie Bit-Vektor-Verbund (*hash filter join*)

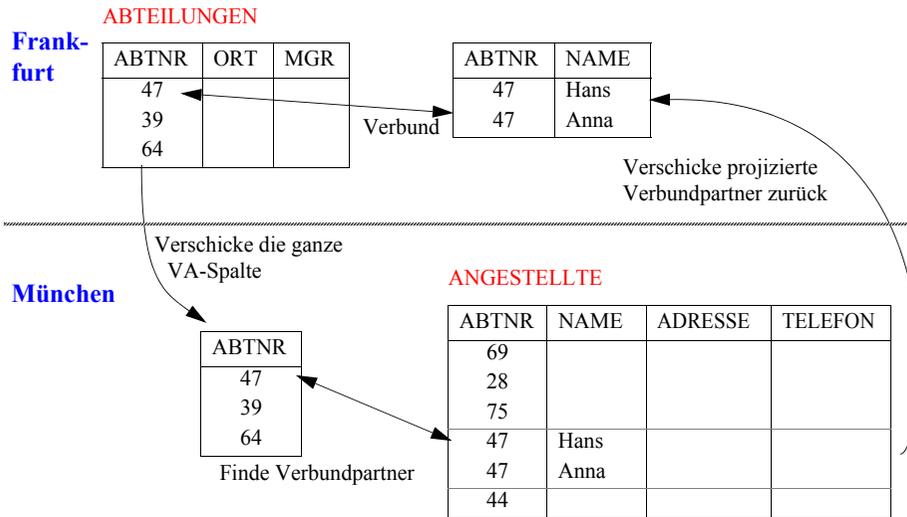
• Semi-Verbund

- Versenden einer Liste der VA von R zum Knoten S
- Ermitteln der Verbundpartner in S und Zurückschicken zum Knoten von R
- Durchführung des Verbundes

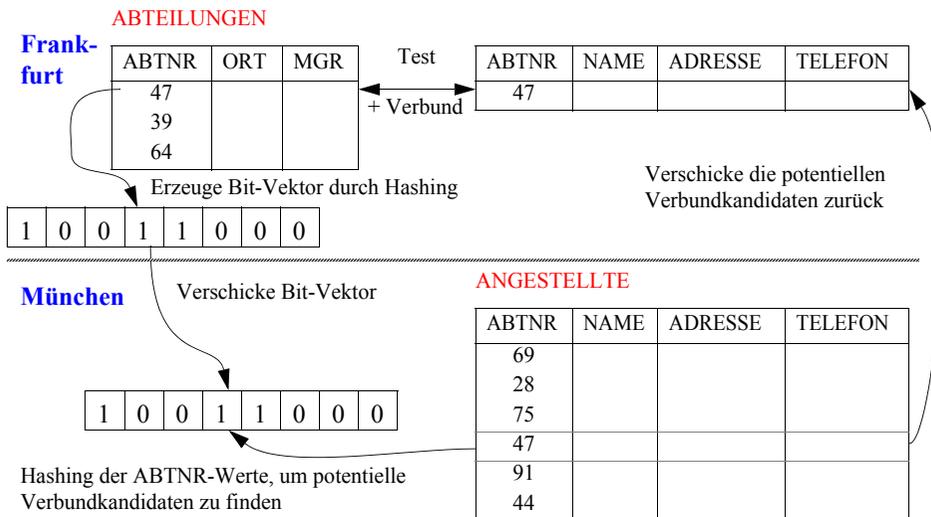
• Bit-Vektor-Verbund

- ähnlich wie Semi-Verbund, nur Versenden eines durch Hash-Funktion erstellten Bitvektors (Bloom-Filter)
- Rücksenden einer Obermenge der Verbundpartner in S

Semi-Verbund

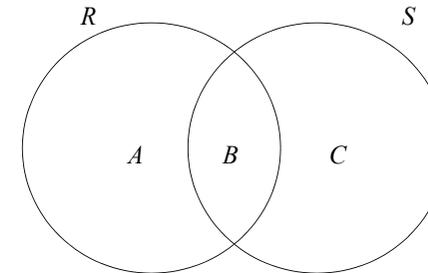


Bit-Vektor-Verbund



Mengenoperationen³

- Welche Mengenoperationen werden benötigt?



R, S vereinigungsverträgliche Eingabeströme
 A, B, C Elementmengen

Operationsergebnis	Übereinstimmung in allen Attributen	Übereinstimmung in einem oder mehreren Attributen
A	Differenz (R-S)	Anti-Semiverbund (S, R)
B	Durchschnitt	Verbund, Semiverbund (S, R)
C	Differenz (S-R)	Anti-Semiverbund (R, S)
A, B		linksseitiger Äußerer Verbund
A, C	Anti-Differenz	Anti-Verbund
B, C		rechtsseitiger Äußerer Verbund
A, B, C	Vereinigung	symmetrischer Äußerer Verbund

- Welche Algorithmen lassen sich für diese Mengenoperationen heranziehen?

- Was muß jeweils verglichen werden?
- Wie läßt sich eine Verbindung zu den Verbundalgorithmen herstellen?

³ Graefe, G.: Query evaluation techniques for large databases, ACM Computing Surveys 25:2, 1993, pp. 73-170

Mengenoperationen (2)

• Binäre Matching-Operationen

- erledigen grundsätzlich dieselbe Aufgabe:
„one-to-one matching operations“
- ein Eingabe-Element trägt zur Ausgabe abhängig von seinem „Match“ mit einem anderen Eingabe-Element bei
- Operationen erfordern immer wieder die gleichen Schritte und können deshalb mit denselben Algorithmen implementiert werden
- Mengen- und Verbundoperationen sind eng miteinander verwandt!

• Gleiche logische Vorgehensweise

- aus R und S werden drei Elementmengen gebildet: A, B, C
- Elemente in B passen zueinander!
- Wie können diese drei Elementmengen gebildet werden?
 - mit Schleifeniteration
 - mit Mischmethode
 - mit Hash-Methode

• Vereinheitlichtes Realisierungskonzept

- Vergleich von Verbund- vs. Primärschlüssel-Attributen
- Gemeinsamkeit: Sätze werden auf der Basis von Attributwerten gruppiert
- Dabei sind einige unäre Operationen mit speziellen Maßnahmen möglich
 - Gruppierung und Sortierung erlaubt einfache Duplikateliminierung
 - Bei Aggregation wird ein Attributwert pro Gruppe bestimmt
 - Beim Verbund ist die Gruppierung der potentiellen Verbundpartner kosteneffektiv (entweder in Partitionen oder einer Sortierordnung)
 - Bei Mengenoperationen können die Elementmengen A, B, C gefunden werden; dabei wird Duplikateliminierung möglich

Zusammenfassung

• Selektionsoperationen

- vorhandene Zugriffspfadtypen erfordern zugeschnittene Operationen und effiziente Abbildung
- Kombination verschiedener Zugriffspfade möglich (TID-Algorithmus)

• Allgemeine Klassen von Auswertungsverfahren für binäre Operationen

- **Schleifeniteration** (nested iteration)
- **Mischmethode** (merge method)
- **Hash-Methode** (hashing)

• Viele Optionen zur Durchführung von Verbundoperationen

- Nested-Loop-Verbund
- Sort-Merge-Verbund
- Hash-Verbund
- und Variationen

• Mengenoperationen

- **prinzipiell Nutzung der gleichen Verfahrensklassen**
- Variation der Vergleichsdurchführung

• Erweiterungsinfrastruktur in objekt-relationalen DBS

- Einbringen von benutzerdefinierten Funktionen und Operatoren
- Verallgemeinerung: benutzerdefinierte Tabellenoperatoren mit n Eingabetabellen und m Ausgabebetabellen