

Relasjonsalgebra.

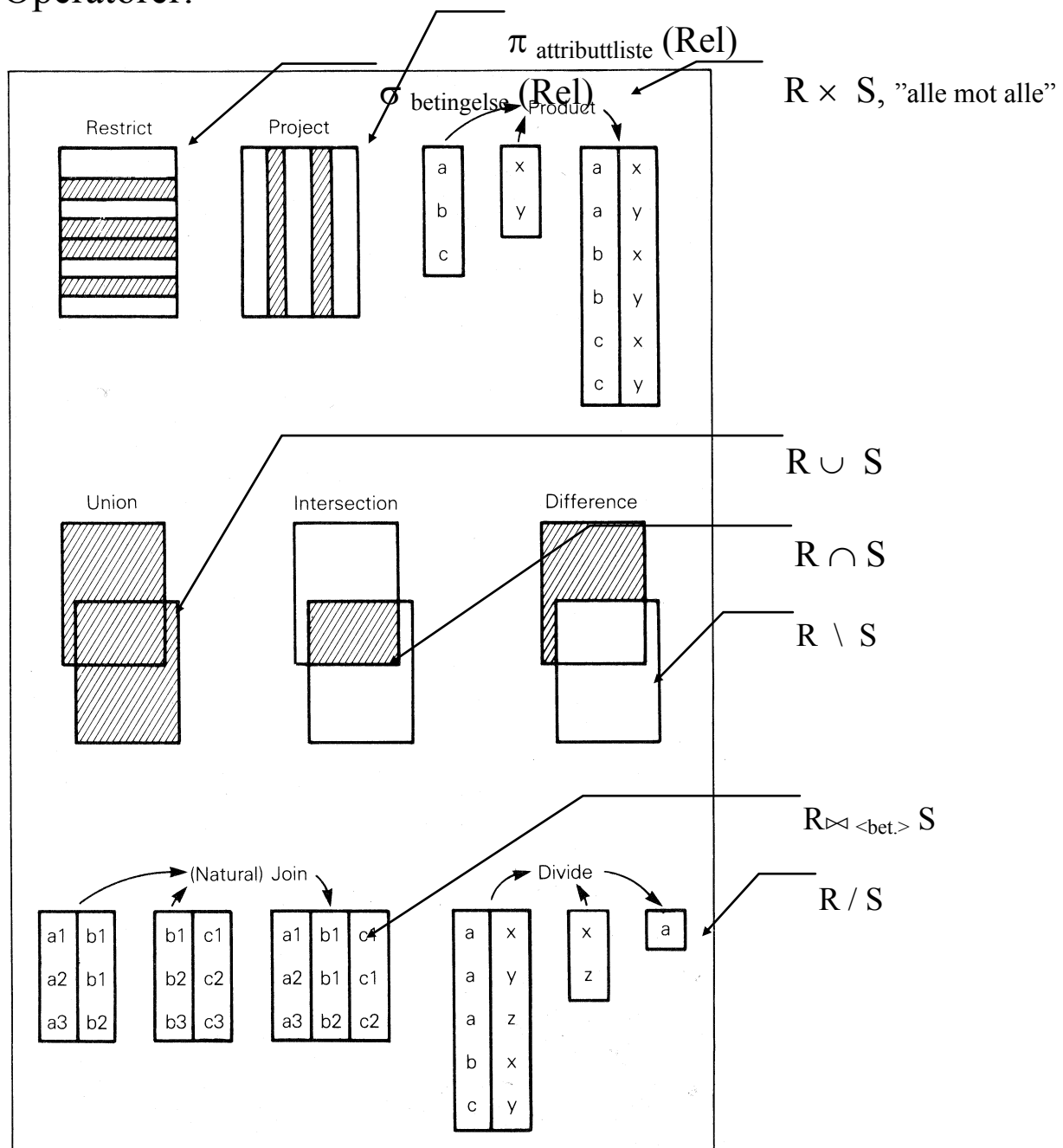
Hva?

Relasjonsalgebra består av et sett med høynivås operatører som kan brukes til å manipulere med relasjoner (slå sammen to tabeller, selektare data etc.). Tankegangen er viktig å kjenne godt til dersom man skal arbeide seriøst med relasjonsdatabaser.

Brukes til bl.a.

- å forstå grunnlaget for relasjonsdatabaser og -spørrespråk
- å forstå hvorledes SQL-setninger kan utføres (og effektiviseres/optimaliseres)
- forstå og arbeide med datavarehus
- forstå og arbeide med distribuerte databaser

Operatorer:



(fra Date: Introduction to Database Systems, utgave 6)

Relasjonsalgebra - vanlige operasjoner.

Mengdeoperasjoner:	<i>Notasjon, variant 1</i>	<i>Notasjon, variant -2</i>
Union	$R \cup S$	R union S
Snitt	$R \cap S$	R intersect S
Mengdedifferanse	$R - S$ $R \setminus S$	R difference S R minus S
Mengdeprodukt, kartesisk produkt ("alle mot alle")	$R \times S$	R product S R times S
<i>Spesielt for relasjoner:</i>		
Horisontalt utvalg (sigma)	$\sigma_{\langle \text{beting.} \rangle}(R)$	R where <bet.> R where <bet.>
Vertikalt utvalg (pi)	$\pi_{\langle \text{feltliste} \rangle}(R)$	R[<feltliste>]
Mengdedivisjon. (Gitt R[c,d] og S[d]. c er med i mengden R dividert med S hvis c i R forekommer sammen med alle d-er som finnes i S.)	$R \div S$ R / S	R divideby S
<i>Spesialiteter av produkt:</i>		
θ -join (produkt med en eller annen betingelse på kompatible attributter, f.eks. >, <, og komb.)	$R \bowtie_{\langle \text{bet.} \rangle} S$	R join<betingelse> S (R join S) where <bet.>
Equi-join (θ -opersjonen er =)	" "	" "
Natural join (Equi-join hvor felles attributt kommer bare en gang) ** den mest vanlige jointypen **	" "	" "
Varianter for produkt:		
Outer join, normalt venstre.. (alle i R, samt alle fra S som oppfyller koblingsbetingelsen)	$R \Join S$	R left join<bet.> S
Full join (alle i R, alle i S, samt alle som oppfyller koblingsbet.)	$R \Join S$	R full join<bet.> S
Semijoin (de i R som tilfredsstiller R join<betingelse> S)	$R \Join_{\langle \text{bet.} \rangle} S$	R semijoin<bet.> S

Legg merke til at operasjonene her er på mengder, slik at evt. dublikater tas bort – tilsvarende select distinct i SQL.

Dersom betingelsen er på primær/fremmednøkkelkomb., droppes ofte <bet>.

Eksempler: - med bruk av σ π \times \bowtie

Gitt tabellene:

Avdeling (Avidnr, Avdnavn, Etasjenr)

Ansatt (Ansnr, Etternavn, Fornavn, Avidnr, Postnr)

Post (Postnr, Poststed)

- Skriv ut alt om ansatte med ansattnr > 100

$\sigma_{\text{ansnr} > 100}(\text{Ansatt})$

- Skriv ut ansattnr, etternavn og fornavn:

$\pi_{\text{Ansnr}, \text{Etternavn}, \text{Fornavn}}(\text{Ansatt})$

- Skriv ut ansattnr og etternavn og fornavn for ansnr > 100:

$\pi_{\text{Ansnr}, \text{Etternavn}, \text{Fornavn}}(\sigma_{\text{ansnr} > 100}(\text{Ansatt}))$

evt. delt opp:

$R1 := \sigma_{\text{ansnr} > 100}(\text{Ansatt})$

$R2 := \pi_{\text{Ansnr}, \text{Etternavn}, \text{Fornavn}}(R1)$

- Skriv ut avdnr, avdnavn og ansnr for alle med postnr < 2000:
f.eks.:

$\pi_{\text{avdnr}, \text{avdnavn}, \text{ansnr}}(\sigma_{\text{postnr} < 2000}(\text{Ansatt} \bowtie \text{Avdeling}))$

NB! Kan gjøres på flere andre måter.

- Skriv ut avdnr, avdnavn for avdelinger som holder til i 3. etasje og har minst en ansatt som bor i 'Bø'.
f.eks.:

$R1 := \sigma_{\text{poststed} = \text{'Bø'}}(\text{Post})$

$R2 := (R1 \bowtie \text{Ansatt})$

$R3 := \pi_{\text{avdnr}}(R2)$

$R4 := R3 \bowtie \text{Avdeling}$

$R5 := \sigma_{\text{etasjenr} = 3}(R4)$

$R6 := \pi_{\text{avdnr}, \text{avdnavn}}(R5)$

- Kunne vært samlet i en. Prøv!

- Kunne vært gjort på mange måter

NB! Siden det er mengdeoperasjoner, tas dublikater automatisk bort

- Finn avdelingsnr for avdelinger uten ansatte

$\pi_{\text{Avidnr}}(\text{Avdeling}) \setminus \pi_{\text{Avidnr}}(\text{Ansatt})$

- Finn avdelinger uten ansatte – ta med både avdelingsnr, -navn og etasje

$(\pi_{\text{Avidnr}}(\text{Avdeling}) \setminus \pi_{\text{Avidnr}}(\text{Ansatt})) \bowtie \text{Avdeling}$

Eksempler, forts: $\sigma\pi \times$ \bowtie

Gitt tabellene:

Avdeling (Avdnr, Avdnavn, Etasjenr)

Ansatt (Ansnr, Etternavn, Fornavn, Avdnr, Postnr) - fast ansatte

Post (Postnr, Poststed)

MAnsatt (Ansnr, Etternavn, Fornavn, Avdnr, Postnr) - midlertidige.

- Lag en liste over evt. som **både** er midlertidige og faste ansatte
 $\text{MAnsatt} \cap \text{Ansatt}$ - NB! Forutsetter at like Ansnr \implies likt på de andre
- Lag en liste over både midlertidige og faste ansatte
 $\text{MAnsatt} \cup \text{Ansatt}$ - NB! Forutsetter at like Ansnr \implies likt på de andre.
 Evt. dublikater forsvinner.
- Hvis noen er både midlertidige og faste, skal de slettes fra midlertidige
 $\text{MAnsatt} := \text{MAnsatt} \setminus \text{MAnsatt} \cap \text{Ansatt}$ - NB! Forutsetter ... som over.
- Finn alle postnr, og liste over de som har dette postnr

$$\pi_{\text{Postnr}}(\text{Post}) \bowtie \pi_{\text{Postnr, Ansnr, ansnavn}}(\text{Ansatt})$$

- Skriv ut alle ansatte i avdelingen for Gressklippere
enklest:

$$\text{Ansatt} \triangleright \sigma_{\text{avdnavn} = \text{'Gressklippere'}}(\text{Avdeling})$$

- Kombiner alle etternavn med alle fornavn

$$\pi_{\text{fornavn}}(\text{Ansatt}) \times \pi_{\text{etternavn}}(\text{Ansatt})$$

- Finn postnr som er slik at minst en fra hver avdeling har dette postnr

$$\pi_{\text{avdnr, postnr}}(\text{Ansatt}) / \pi_{\text{avdnr}}(\text{Avdeling})$$

- Finn avdelingsnr som er slik at de har ansatte fra alle aktuelle postnr

$$\pi_{\text{avdnr, postnr}}(\text{Ansatt}) / \pi_{\text{postnr}}(\text{Post})$$



Relasjonsalgebra vs. SQL

Det kan vises at alle utsagn som skrives i SQL kan skrives i relasjonsalgebra og motsatt – dvs. språkene er like kraftfulle.

Vi finner igjen relasjonsoperatorene i vanlige SQL-setninger:

Eks. 1:

```
SELECT avdnr, avdnavn, ansnr, ansnavn
FROM avd, ans
WHERE avd.avdnr = ans.avdnr
AND etasjenr = 3;
```

hvilken relasjonsoperator?

hvilken relasjonsoperator?

hvilken relasjonsoperator i
utgangspunktet?

hvilken relasjonsoperator?

Eks. 2:

```
SELECT avdnr
FROM ans
WHERE NOT EXISTS
  (SELECT postnr
   FROM post
   WHERE NOT EXISTS
     (SELECT postnr
      FROM ans ans2
      WHERE ans2.ansnr = ans.ansnr and ans.postnr = post.postnr))
```

dvs. “Selekter ut avdnr som er slik at det ikke eksisterer postnr som er slik at det ikke finnes korresponderende avdnr/postnr”

skriv en setning i relasjonsalgebra som gir samme resultat:

Noen regler/egenskaper i rel. algebra.

(temmelig teoretisk

- **Lukkethet:**

Relasjonsalgebra er lukket under relasjonsoperatorene.

- **Assoiativitet:**

Union og snitt er assosiative:

$$(R \cup S) \cup T \equiv R \cup (S \cup T) \text{ og}$$

$$(R \cap S) \cap T \equiv R \cap (S \cap T).$$

Med vår def. av produkt er også denne assosiativ.

- **Kommutativitet:**

Union og snitt er kommutative:

$$R \cup S \equiv S \cup R \quad \text{og} \quad R \cap S \equiv S \cap R$$

Det samme gjelder produkt og inner join.

Derimot er ikke f.eks. mengdedifferanse og outer join.

- **Distributivitet:**

I matematikk f.eks. $a \cdot (b+c) = a \cdot b + a \cdot c$,

$\sqrt{x \circ y} = \sqrt{x} \circ \sqrt{y}$ er distributiv når operatoren \circ er \cdot eller $/$, men ikke dersom den er $+$ eller $-$.

- Restrict er distributiv mht. join i enkle operasjoner, f.eks.

$$\sigma_{\langle R_betingelse \rangle \text{ and } \langle S_betingelse \rangle} (R \bowtie S) \equiv \sigma_{\langle R_betingelse \rangle} (R) \bowtie \sigma_{\langle S_betingelse \rangle} (S)$$

- Prosjeksjon er distributiv mht. join (bortsett fra attributter som deltar i koblingen),

$$\pi_{\langle R_attrliste \rangle \cup \langle S_attrliste \rangle} (R \bowtie S) \equiv \pi_{\langle R_attrliste \rangle} (R) \bowtie \pi_{\langle S_attrliste \rangle} (S)$$

disse brukes bl.a. i forbindelse med optimaliseringsstrategier

Utvidelser av relasjonsalgebra, renames.

Renames

- Vi har tidligere definert nye relasjoner ved tilordning, f.eks.
 $\text{Minavdeling} := \sigma_{\text{avdelingsnr} = 3}(\text{Avdeling})$.
- Vi kan også lage en kopi av en relasjon, f.eks.
 $S2 := S1$, ligner alias, men lager kopi i stedet for å glemme det gamle navnet.
- Vi kan bruke operatoren ρ (rho, gresk r) for å gi et nytt navn på en beregning eller relasjon.
 - $\rho_{\text{Minavdeling}}(\sigma_{\text{avdelingsnr} = 3}(\text{Avdeling}))$, er det samme som tilordningen over.
 - $\rho_{\text{NyAnsatt}}(\text{Ansatt})$ gjør at vi får en ny relasjon NyAnsatt.
 - Ønsker vi andre attributtnavn, er det også mulig, f.eks.
 $\rho_{\text{Employee}}(\text{Empno}, \text{LastName}, \text{FirstName})(\pi_{\text{Ansattnr}, \text{Etternavn}, \text{Fornavn}}(\text{Ansatt}))$.

Beregninger via renames

Dermed kan man også bruke renames for å gjøre beregninger. Vi tenker oss at vi har en varerelasjon med bl.a. varenr og lengde i cm på varen. Vi vi ha en relasjon som passer for England, og vi skal ha med varenr og varenr og lengde både i cm og i tommer.

- $\rho_{\text{EnglishItem}}(\text{Itemno}, \text{CmLength}, \text{InchLength})(\pi_{\text{Varenr}, \text{Lengde}, \text{Lengde} * 2.54}(\text{Vare}))$.

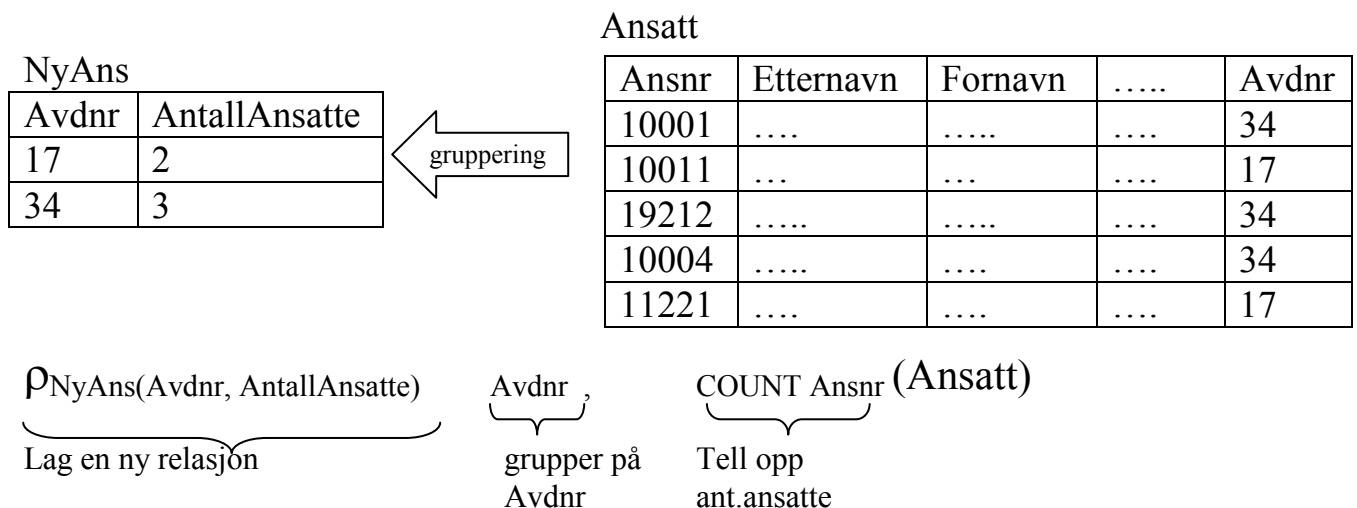
Utvidelser av relasjonsalgebra, gruppering.

Summering / gruppering

Tips: gjør en rask repetisjon på hvordan dette gjøres i SQL.

Siden vi her har bruk for et attributt for å «ta vare på» resultatet f.eks. av en summering i en ny tabell, trenger vi en ny relasjon og nytt/nye attributter. Eksempel:

Tell opp antall ansatte i hver avdeling. Legg den i en ny relasjon med NyAns, med attributtene Avdnr og AntallAnsatte.



Skal vi gjøre flere grupperingsoperasjoner, er det også mulig. Anta at vi har et attributt Lønn for hver ansatt. Vi kan da skrive

$\rho_{\text{NynyAns}}(\text{Avdnr}, \text{AntallAnsatte}, \text{SumLønn}) \text{ Avdnr}, \text{COUNT Ansnr SUM Lønn} (\text{Ansatt})$

Tilsvarende hvis man skal gruppere på flere attributter, f.eks. primært på fylke, og så på kommune innenfor hvert fylke. Kommaen blir da skilletegn mellom attributtet det grupperes på og operatorene.

Hvilke grupperingsoperasjoner? De vanlige:

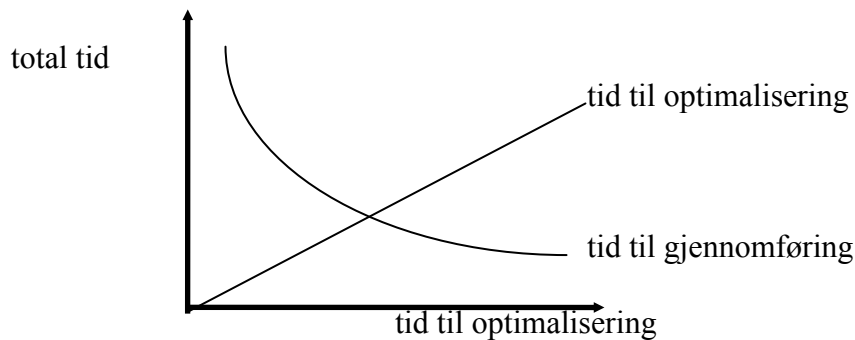
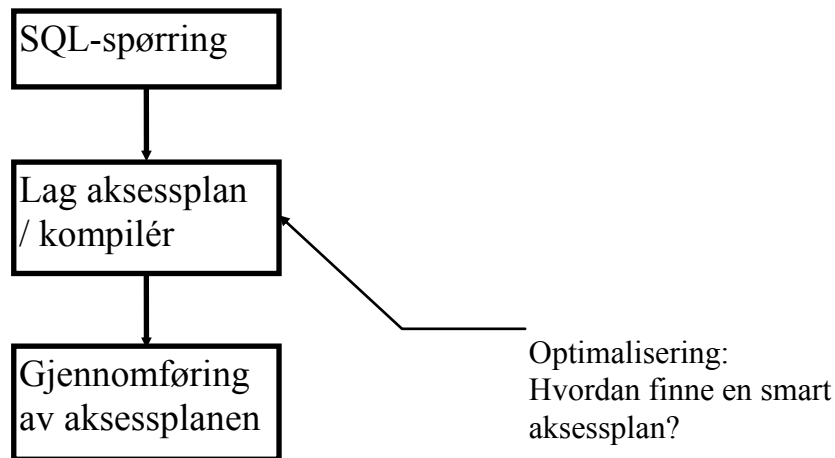
- COUNT
- SUM
- AVG
- MIN
- MAX

Bruksområder for relasjonsalgebra

- allmenn begrepsdannelse
- forståelse av hva som foregår "bak kulissene"
- mulighet for å sjekke kompletthet av et språk
- forståelse av optimaliseringsprosesser
- forståelse av distribuerte databaser (bl.a. hvordan distribuere data i et nettverk)
- kunne vært brukt som spørrespråk, inkl. endring av databasen, men neppe særlig praktisk.
- m.m.

Optimalisering

Hvorledes gjennomføres en spørring?



====> Kan ikke analysere alt for mye for å finne den optimale aksessplanen !!

Spørringer kan være

- adhoc, en kompilering + direkte gjennomføring
- lagrede, en kompilering + $\sum_0^m \{ \text{gjennomføring} \}$

Tidkrevende / "tunge" operasjoner:

- sortering
- fjerning av dublikater
- gruppering
 - oftest lønnsomt å sortere først
- kobling (join)
 - "rå kraft" (alle mot alle, $O(n^2)$)
 - sorter + sammenlign (evt. via fletting)
 - sammenlign via indekser (full søk på en relasjon, via indeks på den andre)
 - hash-metodikk

Optimalisering - hvorledes?

Skritt i optimaliseringsprosessen:

- **nedbrytning** til enkeltoperasjoner og kanonisk form.
- **analyse** med omformulering i hht. lovlige regler, bl.a. ut fra
 - betingelser i where-setningen (= vs. > etc.)
 - hvilke attributter som skal med i resultatet
 - krav til sortering og gruppering
 - antall tupler pr. relasjon
 - antall bytes i attributter som skal kobles
 - indekser som finnes (og om clustede, unike, ikke-unike)
 - ideelt sett hvor dataene finnes (lokalt, på nettet etc.)
 - ideelt sett også variasjonsbredden i dataene.
- **sammenstilling** til en samlet aksessplan

Enkle optimaliseringskriterier:

- **Overordnet: begrenns data mest mulig før tunge operasjoner utføres.**
- Prosjeksjoner gjøres først, men reelt bare hvis relasjoner trengs til mellomregninger, ellers bare ved å vite hvilke attributter som skal med
- Restrict / select-setninger med = verdi/verdisett gjøres først, og de som begrenser mest gjøres først.

$$\text{Eks: } \sigma_{R.a = \langle v \rangle} (R \bowtie S) \implies \sigma_{R.a = \langle v \rangle} (R) \bowtie S$$

- Operasjoner som bør benytte seg av indekser gjøres først
- Operasjoner som gjelder primærnøkkelen gjøres først
- Sjekk om data i en subselect er stabile eller om de må beregnes for hver gang utenforliggende løkke beregnes.
- Boolske shortcut og forenklinger, f.eks.: $P \text{ and } P \equiv P$; $P \text{ and false} \equiv \text{false}$; $P \text{ or false} \equiv P$; $P \text{ and not } P \equiv \text{false}$; $P \text{ and } (P \text{ or } Q) \equiv P$

NB! Antall kombinasjoner øker svært fort, så effektive avskjæringsmekanismer er nødvendig.

Optimalisering - eksempel:

Gitt en Avdelings- og ansattabell.

Avdeling (avdnr, avdnavn,) 5 attributter, 100 tupler

Ansatt (ansnr, ansnavn, avdnr, adresse, postnr) 15 attributter, 5000 tupler

Finn ansnr & -navn for de(n) i Oste-avd. som har bor på postnr 1855.

```
SELECT  ansnr, ansnavn
FROM    Avdeling, Ansatt
WHERE   Avdeling.avdnr = Ansatt.avdnr
        AND postnr = '1850'
        AND avdnavn = 'Ost'
```

Alternativ 1:

$S1 := \sigma_{\text{Avdeling.Avdnr} = \text{Ansatt.Avdnr}} (\text{Avdeling} \bowtie \text{Ansatt})$
 -- 19 attributter, 500.000 poster kobles, resultat ≤ 5000 poster

$S2 := \sigma_{\text{Postnr} = '1850'} (S1)$

$S3 := \sigma_{\text{Avdnavn} = 'Ost'} (S2)$

$S4 := \pi_{\text{ansnr}, \text{ansnavn}}(S3)$ -- fra 19 til 2 attributter

Alternativ 2:

Skritt 1: a) Sett betingelser som gjelder bare en tabell sammen med
 denne tabellen (kalles detachement = frakobling)

b) Fjern overflødige attributter

Skritt 2: Gjør nødvendige koblinger

Skritt 3: Fjern evt. overflødige attributter

$S1 := \sigma_{\text{Avdnavn} = 'Ost'} (\text{Avdeling})$

-- antagelig bare 1 post

$S2 := \pi_{\text{avdnr}}(S1)$

-- og en kolonne / attributt

$S3 := \sigma_{\text{Postnr} = '1850'} (\text{Ansatt})$

-- en/noen få rader / tupler

$S4 := \pi_{\text{avdnr}, \text{ansnr}, \text{ansnavn}} (S3)$

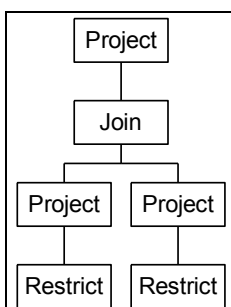
-- og 3 kolonner / attributter

$S5 := S2 \bowtie S4$

-- 1 x 1 / svært få å koble

$S6 := \pi_{\text{ansnr}, \text{ansnavn}} (S5)$

-- kobl.attributtet skal ikke med i resultatet



I en jafs:

$\pi_{\text{ansnr}, \text{ansnavn}} ($

$\pi_{\text{avdnr}} (\sigma_{\text{Avdnavn} = 'Ost'} (\text{Avdeling}))$

\bowtie

$\pi_{\text{avdnr}, \text{ansnr}, \text{ansnavn}} (\sigma_{\text{Postnr} = '1850'} (\text{Ansatt})))$

Optimalisering i databasesystemer:

I teorien:

- Alle spørringer som er ekvivalente burde gi opphav til samme aksessplan.

I praksis:

- Stor forskjell på databasesystemer m.h.t. optimalisering.
- I mange systemer kan man titte på, evt. også endre aksessplanen.
- Med mange prosessormaskiner eller maskiner koblet i klynge blir paralleliserbarhet et viktig kriterium for
 - hele aksessplanen
 - enkeltalgoritmer f.eks. for sortering og kobling

Relasjonskalkyle I.

I stedet for å basere et språk på mengder, kan det baseres på vanlig predikatlogikk. For én variabel:

$\{x \mid P(x)\}$, P er et utsagn som er sant eller usant avhengig av x .
Mer kompliserte uttrykk kan settes sammen av:

	symbol	alternativt symbol
og	\wedge	AND
eller	\vee	OR
negasjon / ikke	\neg eller \sim	NOT
for alle	\forall	FORALL
det eksisterer minst en	\exists	EXISTS

Eksempler:

Gitt: Range of x_{ans} is Ansatt; Range of x_{avd} is Avdeling;

1. Finn ansattnr med $\text{postnr} = 1500$ og $\text{avdnr} = 17$:

$$\{x_{\text{ans.ansattnr}} \mid x_{\text{ans.postnr}} = 1500 \wedge x_{\text{ans.avdnr}} = 17\}$$

2. Finn alle avdelinger som har minst en ansatt:

$$\{x_{\text{avd}} \mid \exists (x_{\text{ans}} \mid x_{\text{avd.avdnr}} = x_{\text{ans.avdnr}})\}$$

3. Finn alle avdelinger uten ansette

$$\{x_{\text{avd}} \mid \neg \exists (x_{\text{ans}} \mid x_{\text{avd.avdnr}} = x_{\text{ans.avdnr}})\}$$

4. Finn om det er én avdeling som har alle ansatte:

$$\{x_{\text{avd}} \mid \forall (x_{\text{ans}} \mid x_{\text{avd.avdnr}} = x_{\text{ans.avdnr}})\}$$

5. Skriv ut avdelinger som har samme navn.

Vi trenger to avdelingsvariable, slik at vi kan sammenligne alle mot alle.

Range of y_{avd} is avdeling;

$$\{x_{\text{avd}} \mid \exists (y_{\text{avd}} \mid (x_{\text{avd.avdnavn}} = y_{\text{avd.avdnavn}}) \wedge (x_{\text{avd.avdnr}} \neq y_{\text{avd.avdnr}}))\}$$

Relasjonskalkyle II.

WFF, frie og bundne variable.

Et syntaktisk og semantisk korrekt utsagn kalles en wff, well-formed formula.

Dersom F er en WFF, så gjelder

$\{\forall x : F(x)\} \equiv \{\neg \exists x \mid \neg F(x)\}$ - dvs.: \forall er egentlig ikke nødvendig. SQL har ikke noen FORALL, derfor må denne omformuleringen brukes.

Tilsvarende :

$\{\exists x \mid F(x)\} \equiv \{\neg(\forall x : \neg F(x))\}$

$\{\forall x : F(x) \wedge G(x)\} \equiv \{\neg \exists x (\neg F(x) \vee \neg G(x))\}$ og

Variablene er av to typer:

Frie variable "løper over alle verdier"

Bundne variable frie variable som knytter seg til formler
som har \forall eller \exists

NB! Sammenlign med frie og bundne kontroller i 4.gen.syst.

Relasjonskalkyle - konkluderende punkter:

- like kraftfullt som relasjonsalgebra
- men helt non-prosydurelt
- SQL bygger mye på tenkningen i relasjonskalkyle.