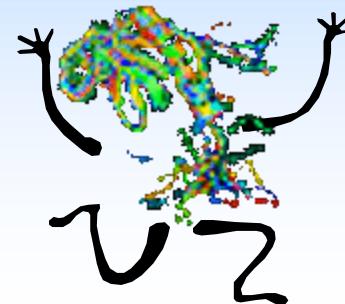


GET THAT PROTEIN!



Eller

TDDE49

Databaser och informationssäkerhet
för bioinformatik

<http://www.ida.liu.se/~TDDE49>

Lärare

- Examinator: Olaf Hartig
- FÖ: Olaf Hartig, Ulf Kargén, Patrick Lambrix
- LA: Adriana Concha, Ulf Kargén
- projekt: Patrick Lambrix, Adriana Concha
- databasadministration: Adriana Concha
- studierektor: Patrick Lambrix
- Olaf/ Adriana på engelska

Kurslitteratur

- Elmasri, R. and Navathe, S. B. Fundamentals of Database Systems, Addison Wesley.
- (Padron-McCarthy, webbkurs på Svenska)
- van Oorschot P: Computer Security and the Internet.
- Anderson R, Security Engineering.
- Adam Shostack: Threat Modeling.
- Lab + projekt: på hemsidan

Databaser

- Ett (av flera) sätt att lagra data i elektronisk format
- Används i det vardagliga livet: bank, bokning av hotell eller resa, sökning i biblioteket, handla
- nyare tillämpningar: multimediadatabaser, geografiska informationssystem, realtiddatabaser

Databaser

- databashanteringssystem (DBMS): en uppsättning program som tillåter en användare att skapa och underhålla en databas
- databassystem = databas + databashanteringssystem

Bioinformatik

- Kända sekvenser samlas i en stor databas. Insamlande och studier av sekvenser och jämförelser av sekvensernas uppbyggnad i olika organismer kallas bioinformatik. Forskningen inom bioinformatik är beroende av avancerad datalogi och matematik. (forskningsrådens strategidokument 2000)

Bioinformatik

- Bioinformatics: research, development, or application of computational tools and approaches for expanding the use of biological, medical, behavioral or health data, including those to acquire, store, organize, archive, analyze or visualize data. (National Institutes of Health)

Bioinformatik

Ämnen på ISMB:

- protein structure and modeling
- sequence motifs, alignments and families
- networks and modeling
- gene structure, regulation and modeling
- sequence and phylogeny
- databases, information and knowledge management

TDDE49 Databaser och informationssäkerhet för Bioinformatik

- Denna kurs: fokus på biologiska databanker

Relation med andra kurser inom TB-programmet:

- förkunskaper: molekylärbiologi,
programmering
- bioinformatik

Årets ändringar i kursen

Inga större ändringar

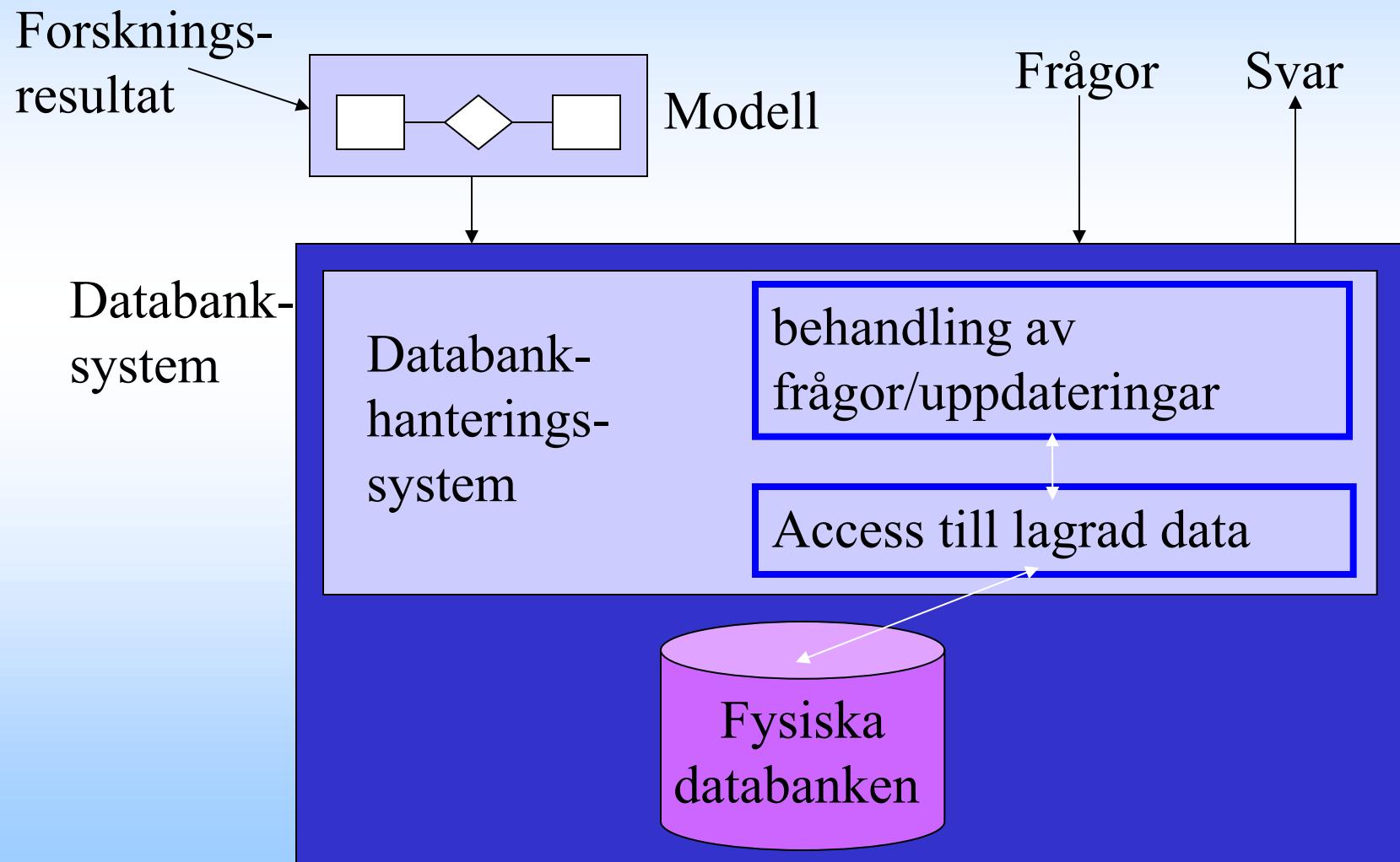
Tidigare ändringar senaste åren:

- Namnbyte
- Informationssäkerhet tillkommer
- Information retrieval, semi-strukturerad data, teoretiska delar av relationsdatabaser utgår.
- Nya labbar för databasdelen.
- Flipped classroom för databasföreläsningarna.

Biologiska databanker

- biologisk data i elektronisk format
- exempel: SWISS-PROT/UniProt, EMBL, DDBJ, PDB, GENBANK, KEGG, ACEDB
- används dagligen i forskningen

Biologiska databanker



Frågeställningar

- Vilken information lagrar man?
- Hur lagras informationen? (hög och låg nivå)
- Hur accessar man informationen?
(användarnivå, systemnivå)
- Hur återställer man en databank efter crash?
- Hur kan flera användare accessa och
uppdatera informationen samtidigt?

Personer

- databankadministratör
 - databankdesigner
 - användare ('end user')
 - programmerare av tillämpningar
-
- DBMS designer
 - utvecklare av verktyg
 - operator, underhåll

1 tgctacccgc gcccgggctt ctggggtgtt ccccaaccac ggcccagccc tgccacaccc
61 cccgcccccg gcctccgcag ctcggcatgg gcgcgggggt gctcgccctg ggcgcctccg
121 agcccgtaa cctgtcgctg gcccacccgc tccccgacgg cgccggccacc gcccgcggc
181 tgctggtgcc cgctgcggcg cccgccttgt tgctgcetcc cgccagcgaa agcccgagc
241 cgctgtctca gcagtggaca ggggcatgg gtctgtat ggcgtctatc gtgtgtctca
301 tcgtggcggg caatgtgctg gtgatgtgg ccatacgccaa gacgcccggg ctgeagacgc
361 tcaccaacct ttcatcatg tccctggcca gcccgcacct ggtcatgggg ctgtgttgg
421 tgccgttcgg ggccaccatc gtgggtgtgg gcccgtggga gtacggctcc ttctctgcg
481 agctgtggac ctcaagtggac gtgtgtgcg tgacggccag catcgagacc ctgtgtgtca
541 ttgccctgga cccgtaccc gccatcacct cgccttcgg ctaccagagc ctgtgtacgc
601 gcccgcgggc gccccccctc gtgtgcaccc tggtggccat ctggccctg tggtccctcc
661 tgcccatctt catgcactgg tggcgccgg agagcgacga ggcgcgcgc tgctacaacg
721 accccaagtg ctgcgacttc gtacccaacc gggcctacgc catgcctcg tccgttagtct
781 cttctacgt gcccgtgtc atcatggct tcgtgtaccc gccccgttcc cgcgaggccc
841 agaagcaggtaa gaagaagatc gacagctgcg agcgcgttt ctcggcgcc ccagegcggc
901 cgcctcgcc ctgccttcg cccgtccccg cgcctcgcc gcccggccga cccccggcc
961 cccgcgcgc cccgcgcacc gcccgcgtt ccaacggggc tgccggtaag cggcgccct
1021 cgcgcctcggtt ggcctacgc gagcagaagg cgctcaagac gtcggccatc atcatggcg
1081 tttcacgtt ctgcgtgtcc tggccaacgt ggtgaaggcc ttccaccgcg
1141 agctggtgcc cggcccttc ttcgttttcaactggctt gggctacgc aactcgccct
1201 tcaacccat catctactgc cgcagccccg acttccgcac ggcctccag ggactgtct
1261 gtcgcgcgc cagggtgtcc cgcggcgcc acgcgcacca cggagacccgg cgcgcgcgc
1321 cgggtgtctt gccccggcc ggaccccccgc catgccttcg gcccgcctcg gacgacgacg
1381 acgacgatgt cgtcgccggcc acggccggcc cgcgcctgtt ggagccctgg gcccgtgtca
1441 acggcgccggc ggcggcgac agcgactcga gcctggacga gcccgtccgc cccggcttcg
1501 cctcgaaatc caagggttag ggcggcgcc gggcgccggatc ctcggccac ggttcccaag
1561 gggaaacgagg agatctgtt ttacttaaga ccgatagcag gtgaactcga agcccaaat
1621 ctcgtctga atcatccgag gcaaagagaa aagccacggc ccgtgcaca aaaaggaaag
1681 ttgggaagg gatgggagag tggctgtatgttcccttg ttg

DEFINITION	Homo sapiens adrenergic, beta-1-, receptor
ACCESSION	NM_000684
SOURCE ORGANISM	human
REFERENCE	1
AUTHORS	Frielle, Collins, Daniel, Caron, Lefkowitz, Kobilka
TITLE	Cloning of the cDNA for the human beta 1-adrenergic receptor
REFERENCE	2
AUTHORS	Frielle, Kobilka, Lefkowitz, Caron
TITLE	Human beta 1- and beta 2-adrenergic receptors: structurally and functionally related receptors derived from distinct genes

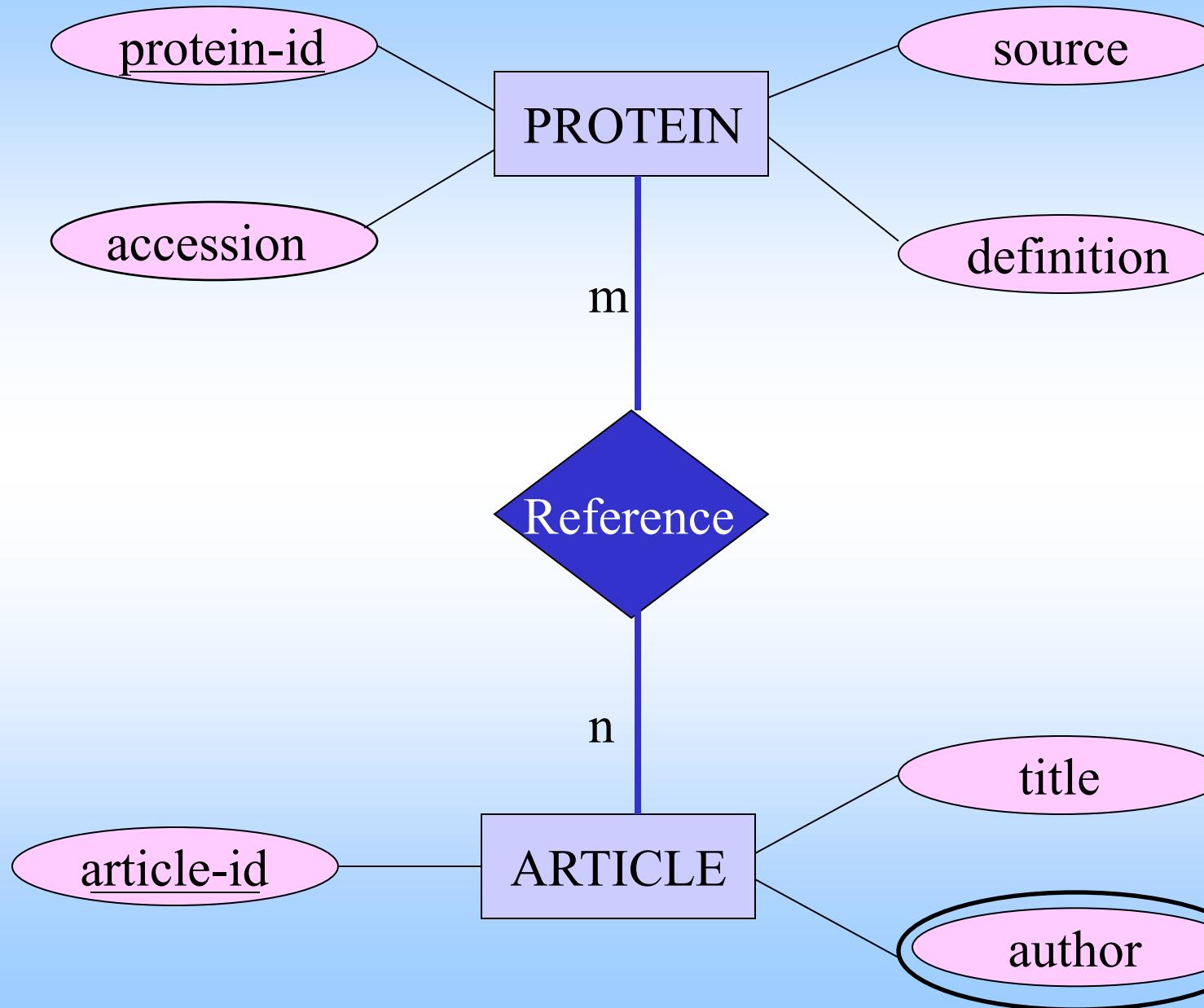
Vilken information lagrar man?

- Modell av verkligheten
 - Entity-Relationship modell (ER)
 - Unified Modeling Language (UML)

Entity-Relationship

- entiteter och attribut
- entitetstyper
- nyckelattribut
- relationer
- kardinalitetsvillkor

Entity-relationship



Hur lagras informationen?

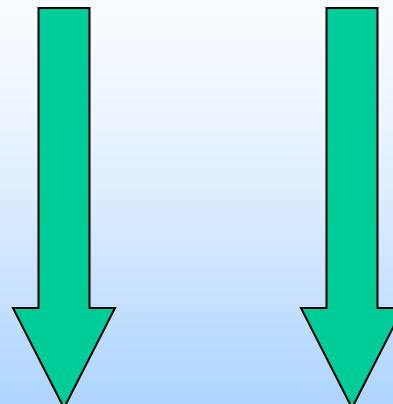
(hög nivå)

Hur accessar man informationen?

(användarnivå)

- Text (IR)
- Semistrukturerad data
- Datamodeller (DB)
- Regler + Fakta (KB)

struktur precision



Text - Information Retrieval

- sökning baseras på ord
- konceptuella modeller:
boolesk, vektor, probabilistisk, ...
- filmodell:
flat fil, inverterad fil, ...

IR - Filmmodell: inverterad fil

inverterad fil

WORD	HITS	LINK
...
adrenergic	32	—
...
cloning	53	—
...
receptor	22	—
...

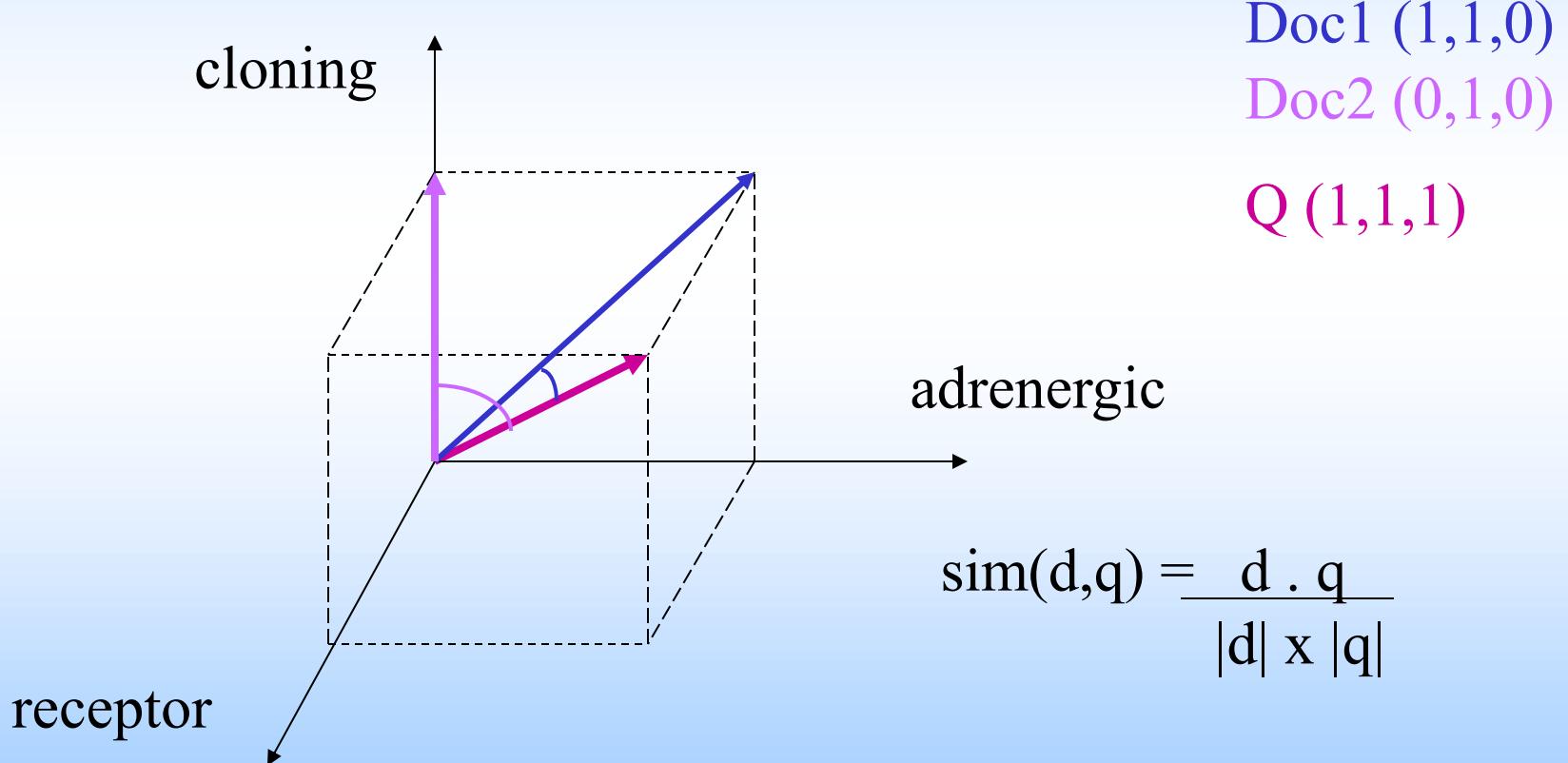
anslagningsfil

DOC#	LINK
...	...
1	—
5	—
...	...
1	—
2	—
5	—
...	...

dokumentfil

DOCUMENTS
Doc1
Doc2
...
...

Vektormodellen (förenklad)



Databaser

- Relationsdatabaser:
 - modell: tabeller + relationsalgebran
 - frågespråk (SQL)
- Objektorienterade databaser:
 - modell: fortlevande objekt,
meddelande, inkapsling, ärvning
 - frågespråk (t.ex. OQL)
- System: GDB (R), ACEDB (OO)

Relationsdatabaser

PROTEIN

PROTEIN-ID	ACCESSION	DEFINITION	SOURCE	REFERENCE	ARTICLE-ID
1	NM_000684	Homo sapiens adrenergic, beta-1-, receptor	human	1 1	1 2

ARTICLE

ARTICLE-ID	AUTHOR	TITLE
1	Frielle	Cloning of the cDNA for the human
1	Collins	Cloning of the cDNA for the human
1	Daniel	Cloning of the cDNA for the human
1	Caron	Cloning of the cDNA for the human
1	Lefkowitz	Cloning of the cDNA for the human
1	Kobilka	Cloning of the cDNA for the human
2	Frielle	Human beta 1- and beta 2-adrenergic receptors
2	Kobilka	Human beta 1- and beta 2-adrenergic receptors
2	Lefkowitz	Human beta 1- and beta 2-adrenergic receptors
2	Caron	Human beta 1- and beta 2-adrenergic receptors

Relationsdatabaser

PROTEIN

PROTEIN-ID	ACCESSION	DEFINITION	SOURCE	REFERENCE	ARTICLE-ID
1	NM_000684	Homo sapiens adrenergic, beta-1-, receptor	human	1 1	1 2

ARTICLE-AUTHOR

ARTICLE-ID	AUTHOR	ARTICLE-ID	TITLE
1	Frielle		Cloning of the cDNA for the human beta 1-adrenergic receptor
1	Collins	1	
1	Daniel		
1	Caron	2	Human beta 1- and beta 2- adrenergic receptors: structurally and functionally related receptors derived from distinct genes
1	Lefkowitz		
1	Kobilka		
2	Frielle		
2	Kobilka		
2	Lefkowitz		
2	Caron		

SQL

```
select source  
from protein  
where accession = NM_000684;
```

Vilka kolumner?
Vilka tabeller?
Vilka rader?

PROTEIN			
PROTEIN-ID	ACCESSION	DEFINITION	SOURCE
1	NM_000684	Homo sapiens adrenergic, beta-1-, receptor	human

SQL

```
select title  
from protein, article-title, reference  
where protein.accession = NM_000684  
and protein.protein-id  
= reference.protein-id  
and reference.article-id  
= article-title.article-id;
```

PROTEIN

PROTEIN-ID	ACCESSION	DEFINITION	SOURCE
1	NM_000684	Homo sapiens adrenergic, beta-1-, receptor	human

REFERENCE

PROTEIN-ID	ARTICLE-ID
1	1
1	2

ARTICLE-TITLE

ARTICLE-ID	TITLE
1	Cloning of the ...
2	Human beta 1- ...

SQL

```
select title  
from protein, article-title, reference  
where protein.accession = NM_000684  
and protein.protein-id  
= reference.protein-id  
and reference.article-id  
= article-title.article-id;
```

PROTEIN

PROTEIN-ID	ACCESSION	DEFINITION	SOURCE
1	NM_000684	Homo sapiens adrenergic, beta-1-, receptor	human

REFERENCE

PROTEIN-ID	ARTICLE-ID
1	1
1	2

ARTICLE-TITLE

ARTICLE-ID	TITLE
1	Cloning of the ...
2	Human beta 1- ...

PROTEIN-ID	ACCESSION	DEFINITION	SOURCE
1	NM_000684	Homo sapiens adrenergic, beta-1-, receptor	human

SQL

```
select title  
from protein, article-title, reference  
where protein.accession = NM_000684  
and protein.protein-id  
      = reference.protein-id  
and reference.article-id  
      = article-title.article-id;
```

PROTEIN

PROTEIN-ID	ACCESSION	DEFINITION	SOURCE
1	NM_000684	Homo sapiens adrenergic, beta-1-, receptor	human

REFERENCE

PROTEIN-ID	ARTICLE-ID
1	1
1	2

ARTICLE-TITLE

ARTICLE-ID	TITLE
1	Cloning of the ...
2	Human beta 1- ...

PROTEIN

PROTEIN-ID	ACCESSION	DEFINITION	SOURCE	REFERENCE	
				PROTEIN-ID	ARTICLE-ID
1	NM_000684	Homo sapiens adrenergic, beta-1-, receptor	human	1	1
1	NM_000684	Homo sapiens adrenergic, beta-1-, receptor	human	1	2

SQL

```
select title  
from protein, article-title, reference  
where protein.accession = NM_000684  
and protein.protein-id  
      = reference.protein-id  
and reference.article-id  
      = article-title.article-id;
```

PROTEIN

PROTEIN-ID	ACCESSION	DEFINITION	SOURCE
1	NM_000684	Homo sapiens adrenergic, beta-1-, receptor	human

REFERENCE

PROTEIN-ID	ARTICLE-ID
1	1
1	2

ARTICLE-TITLE

ARTICLE-ID	TITLE
1	Cloning of the ...
2	Human beta 1- ...

PROTEIN-ID	ACCESSION	DEFINITION	SOURCE	PROTEIN-ID	REFERENCE		ARTICLE-TITLE
					ARTICLE-ID	ARTICLE-ID	
1	NM_000684	Homo sapiens adrenergic, beta-1-, receptor	human	1	1	1	Cloning of the ...
1	NM_000684	Homo sapiens adrenergic, beta-1-, receptor	human	1	2	2	Human beta 1- ...

SQL

```
select title  
from protein, article-title, reference  
where protein.accession = NM_000684  
and protein.protein-id  
      = reference.protein-id  
and reference.article-id  
      = article-title.article-id;
```

PROTEIN

PROTEIN-ID	ACCESSION	DEFINITION	SOURCE
1	NM_000684	Homo sapiens adrenergic, beta-1-, receptor	human

REFERENCE

PROTEIN-ID	ARTICLE-ID
1	1
1	2

ARTICLE-TITLE

ARTICLE-ID	TITLE
1	Cloning of the ...
2	Human beta 1- ...

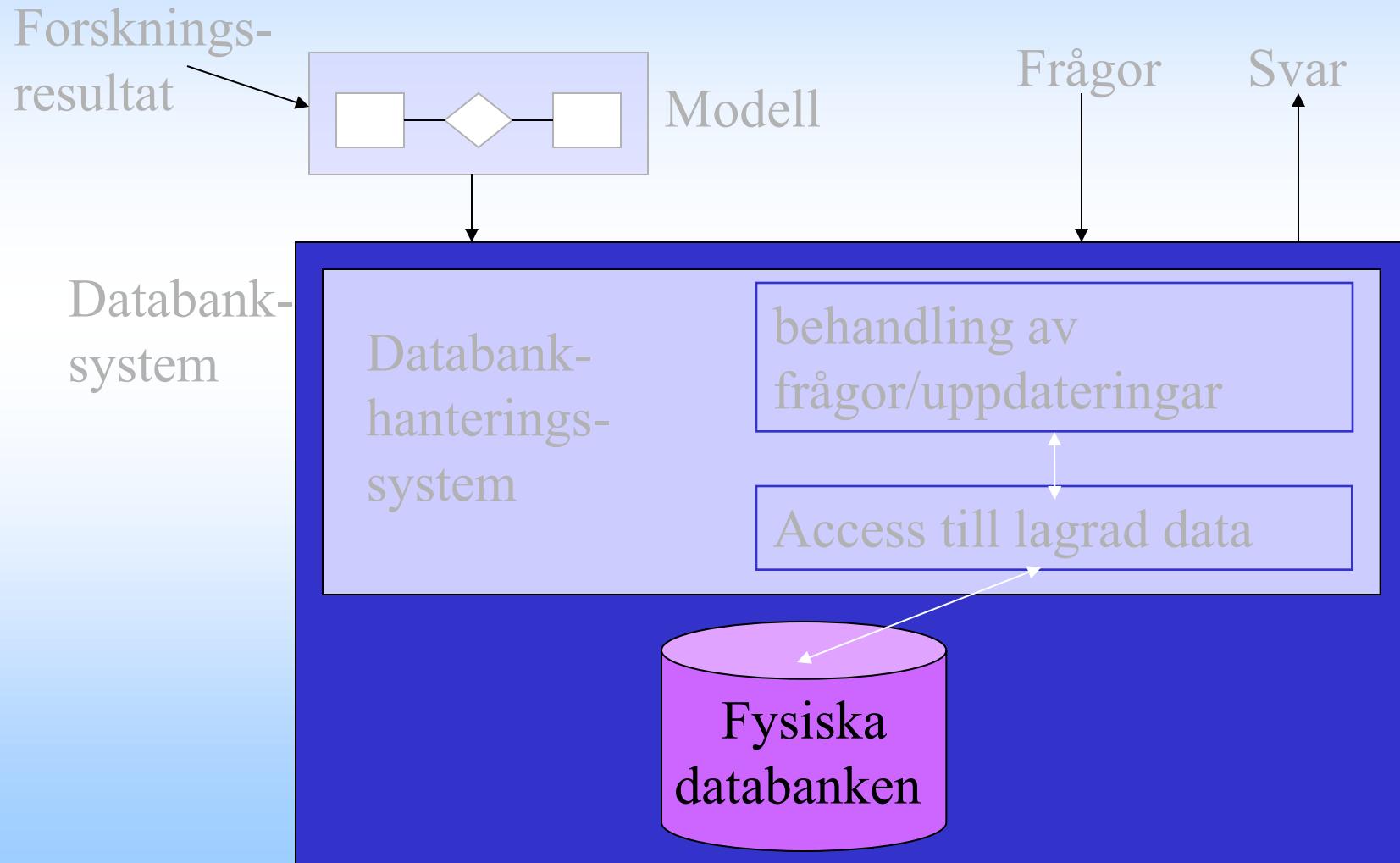
PROTEIN-ID	ACCESSION	DEFINITION	SOURCE	PROTEIN-ID	REFERENCE	ARTICLE-TITLE	TITLE
					ARTICLE-ID	ARTICLE-ID	
1	NM_000684	Homo sapiens adrenergic, beta-1-, receptor	human	1	1	1	Cloning of the ...
1	NM_000684	Homo sapiens adrenergic, beta-1-, receptor	human	1	2	2	Human beta 1- ...

TITLE

Cloning of the ...

Human beta 1- ...

Hur lagras informationen? (läg nivå)



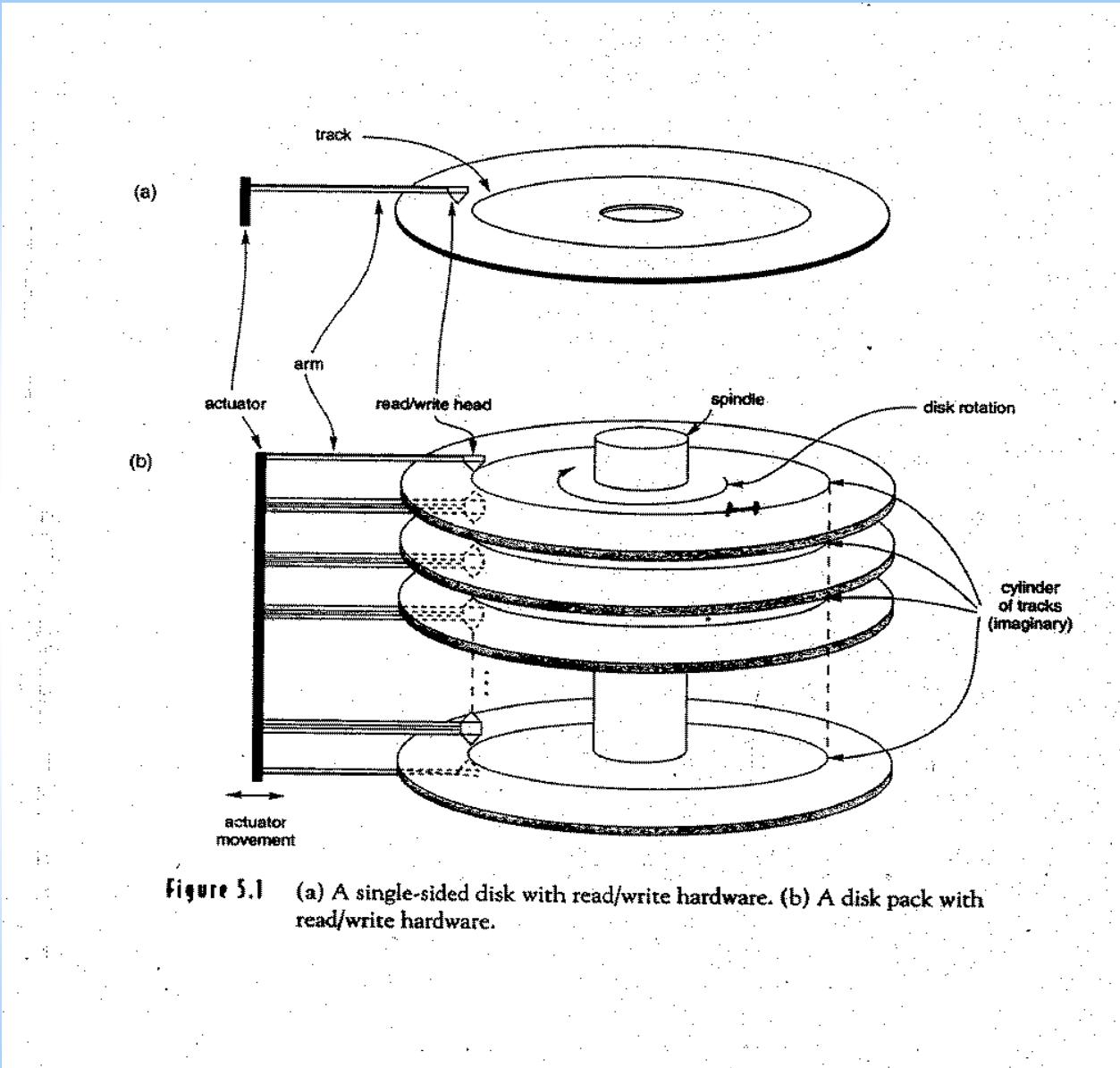
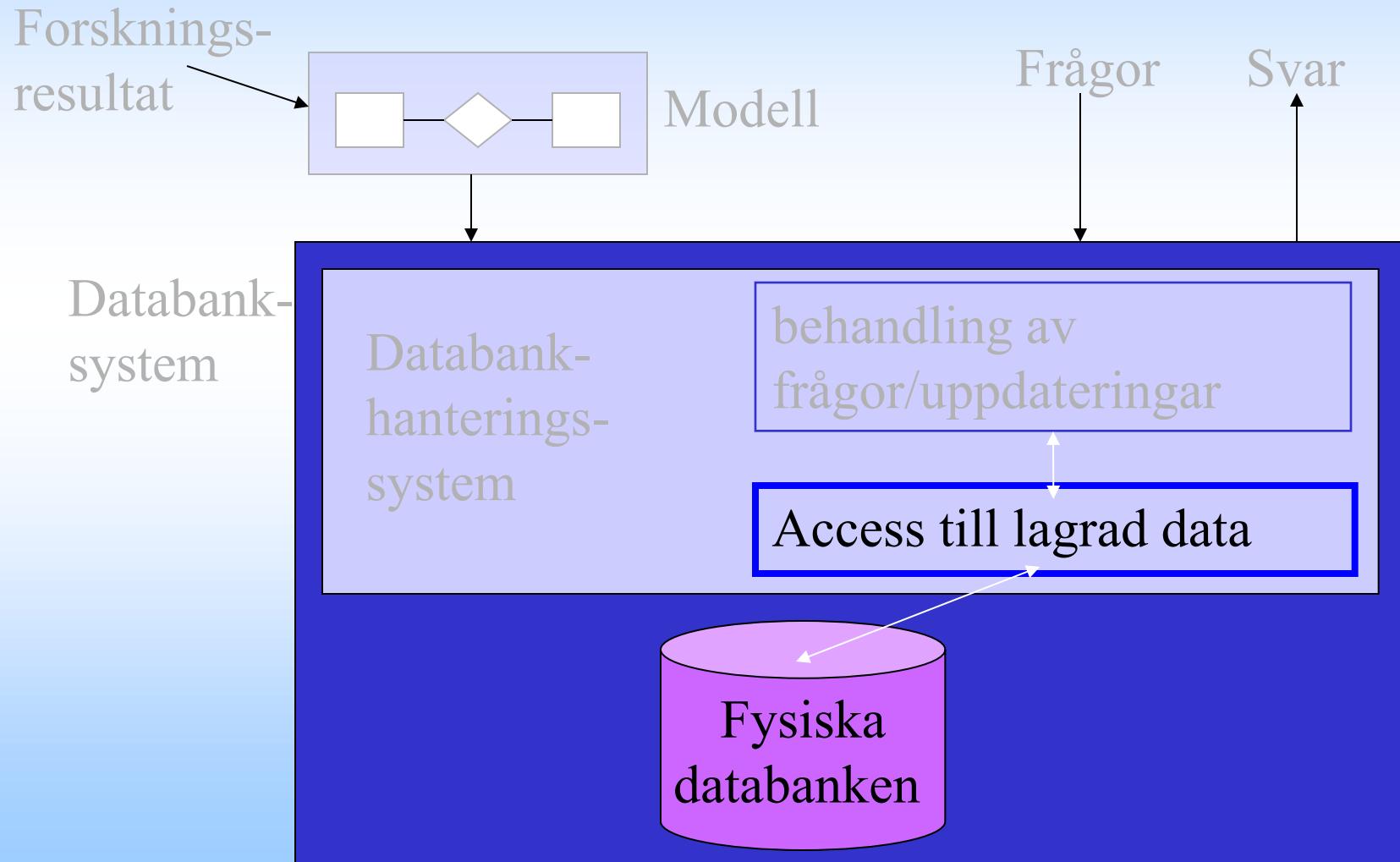


Figure 5.1 (a) A single-sided disk with read/write hardware. (b) A disk pack with read/write hardware.

Hur accessar man informationen? (systemnivå)

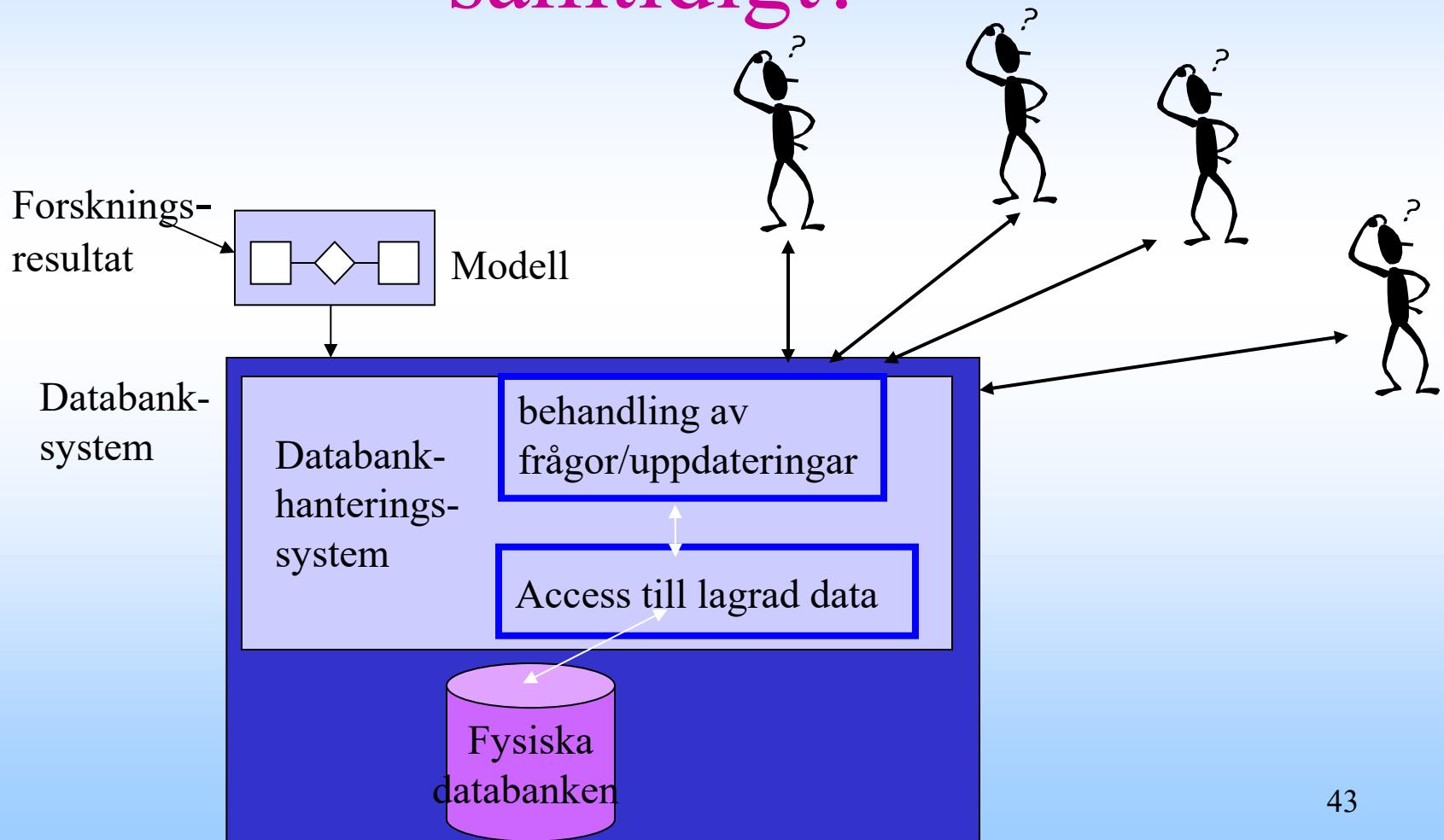


Hur återställer man en databank efter crash?

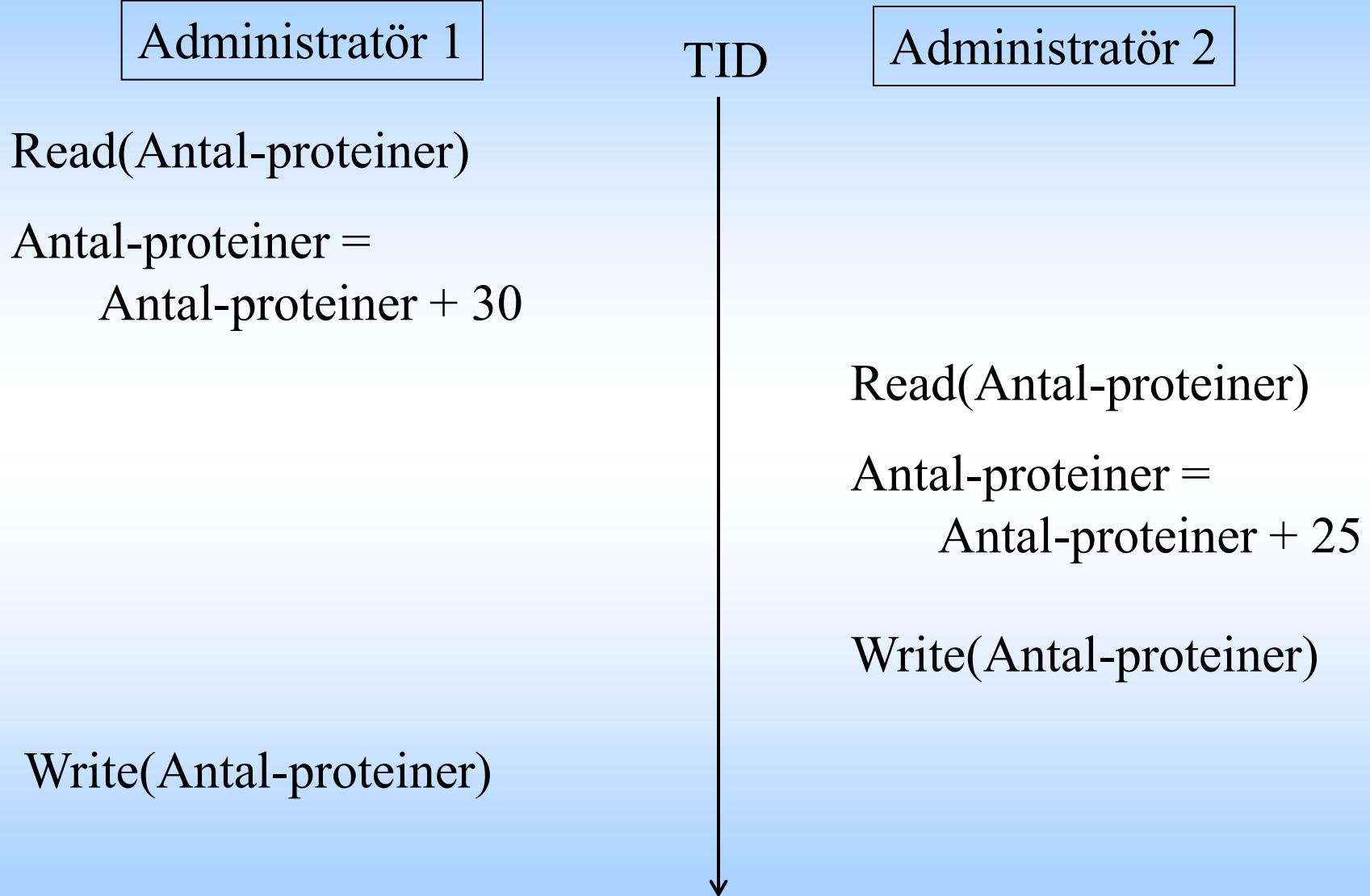
Återställning vid

- datorstop (system crash)
- systemfel
- samtidighetsfel (flera användare)
- skivfel
- katastrofer

Hur kan flera användare accessa och uppdatera informationen samtidigt?



Flera användare



Flera användare

DB



Antal-proteiner: 150

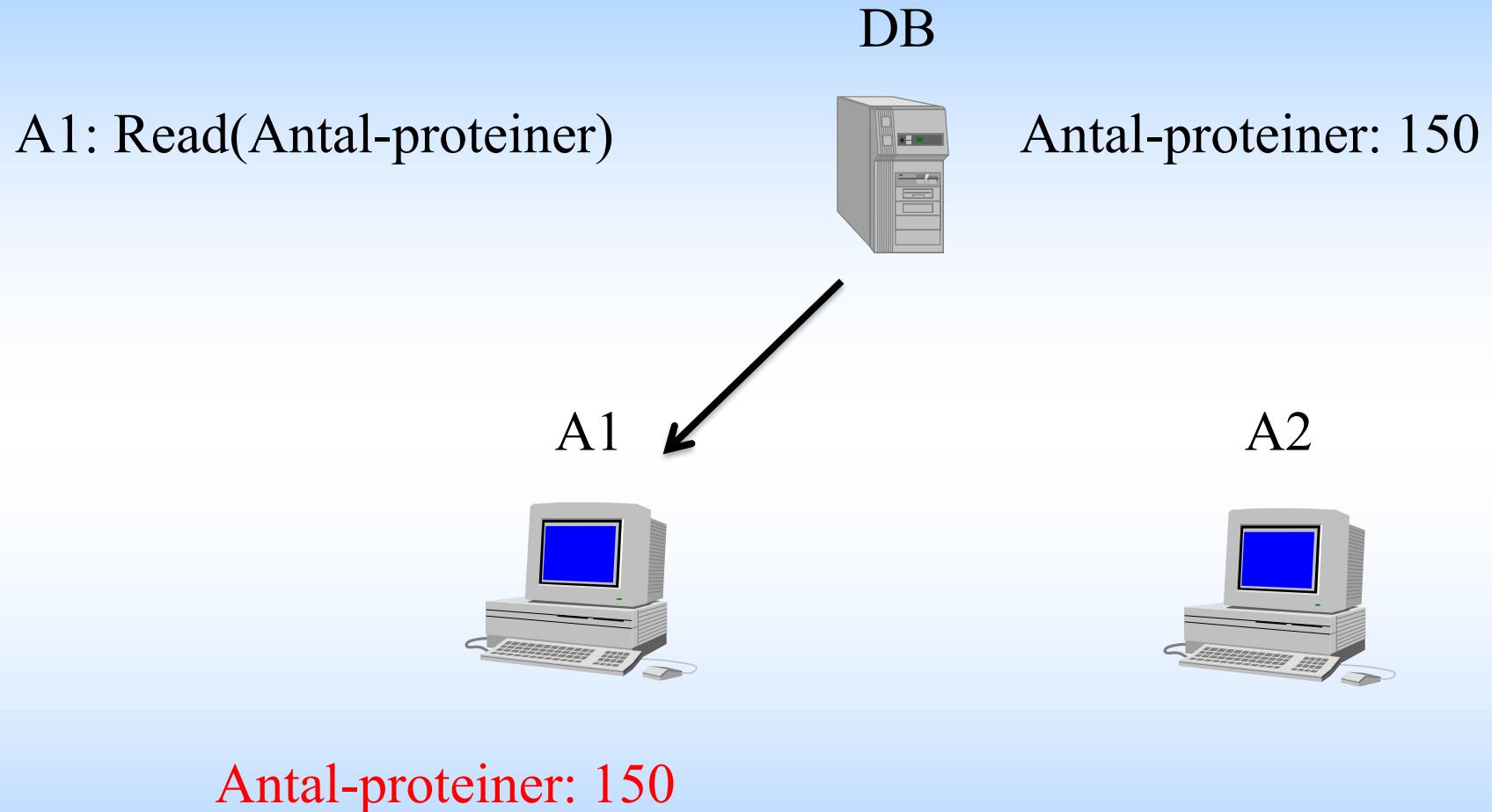
A1



A2



Flera användare



Flera användare

A1: Antal-proteiner =
Antal-proteiner + 30

DB



Antal-proteiner: 150

A1



A2



Antal-proteiner: 150 + 30

Flera användare

A1: Antal-proteiner =
Antal-proteiner + 30

DB



Antal-proteiner: 150

A1



A2



Antal-proteiner: 180

Flera användare

A2: Read(Antal-proteiner)

DB



Antal-proteiner: 150

A1



Antal-proteiner: 180

A2



Antal-proteiner: 150

Flera användare

A2: Antal-proteiner =
Antal-proteiner + 25

DB



Antal-proteiner: 150

A1



Antal-proteiner: 180

A2



Antal-proteiner: 150 +25

Flera användare

A2: Antal-proteiner =
Antal-proteiner + 25

DB



Antal-proteiner: 150

A1



Antal-proteiner: 180

A2



Antal-proteiner: 175

Flera användare

A2: Write(Antal-proteiner)

DB



Antal-proteiner: 150

~~150~~

A1



Antal-proteiner: 180

A2



Antal-proteiner: 175

Flera användare

A2: Write(Antal-proteiner)

DB



Antal-proteiner: 175

A1



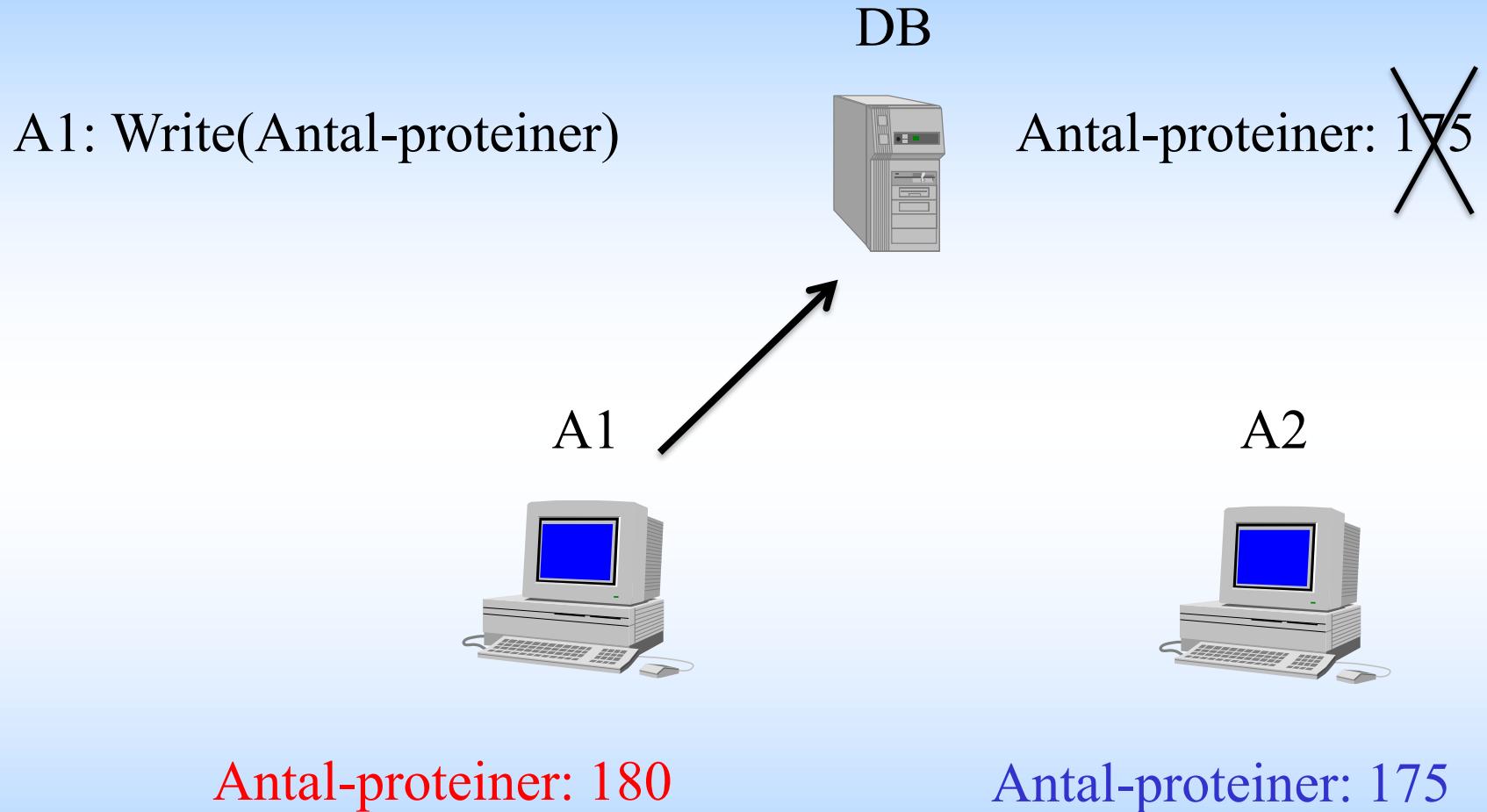
Antal-proteiner: 180

A2

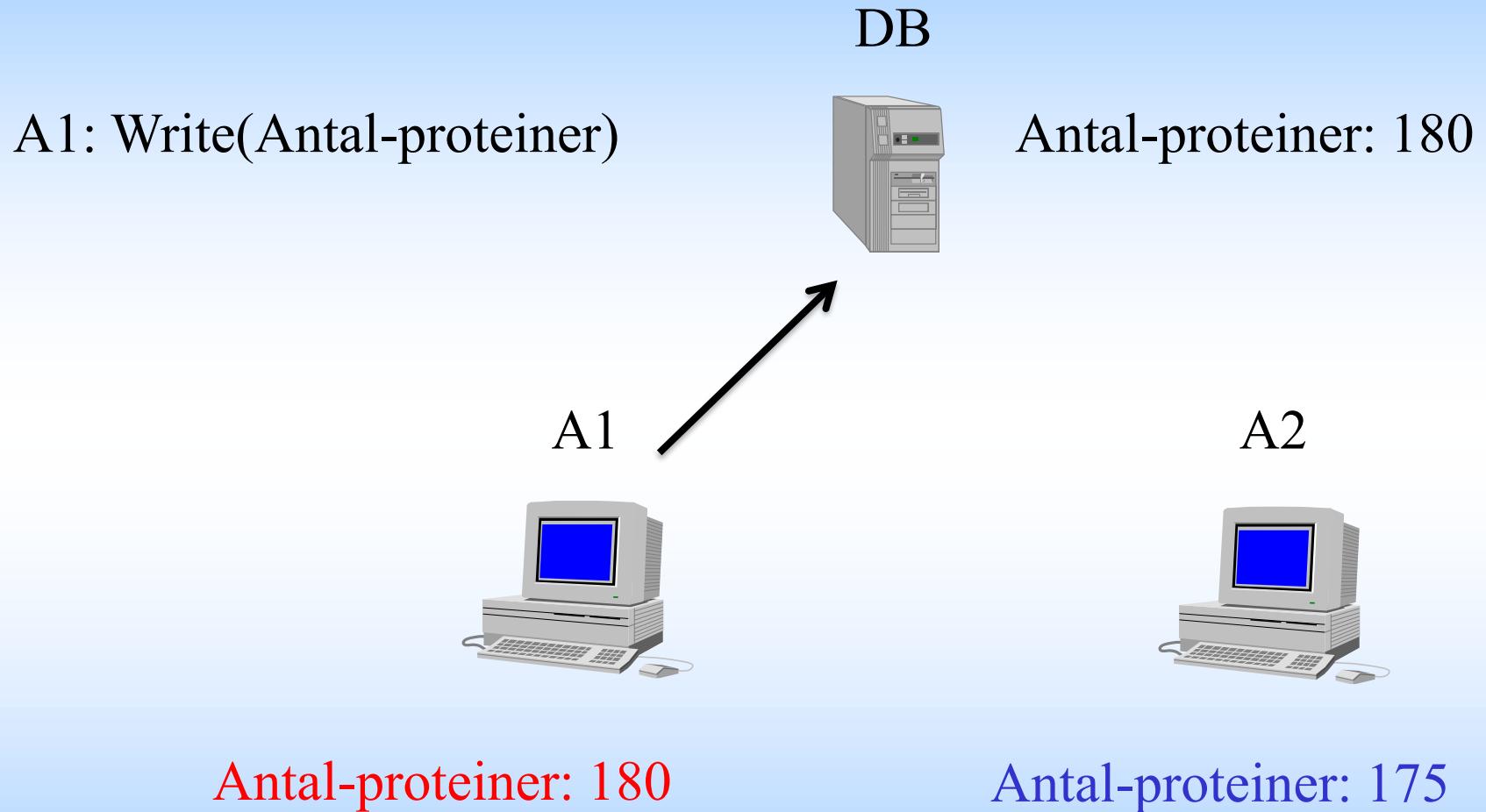


Antal-proteiner: 175

Flera användare



Flera användare



Flera användare

DB



Antal-proteiner: 180

Antal-proteiner: $150 + 30 + 25 = 205$

Informationssäkerhet

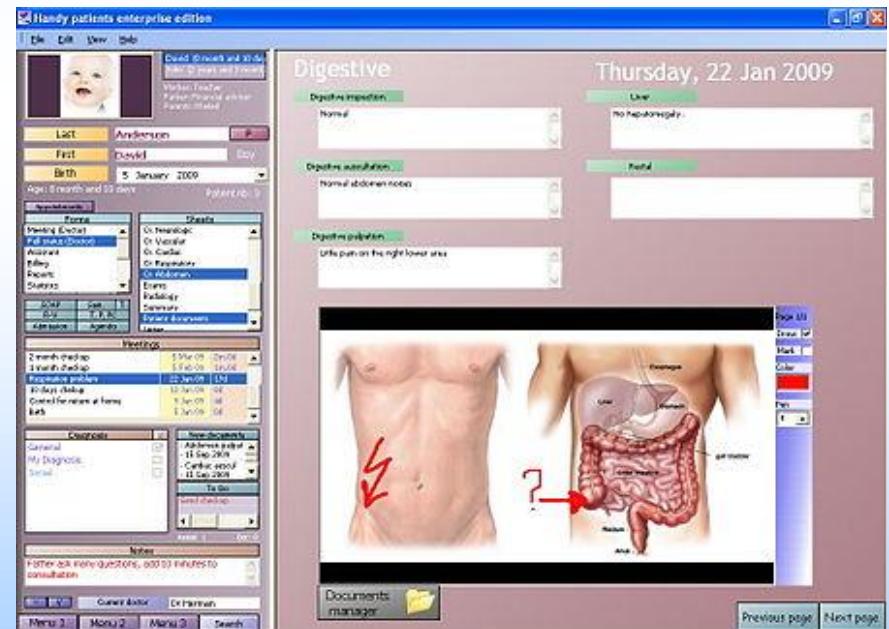
Sekretess

Enbart behöriga användare får ta del av informationen

Integritet

Korrekt och fullständig information

Tillgänglighet



Kursöversikt - FÖ

- Introduktion
- Relationsdatabaser och SQL
- Datamodellering, ER/EER diagram
- Att gå från EER diagram till relationsscheman

Kursöversikt - FÖ

- Informationssäkerhet
- Databassäkerhet

Kursöversikt - LA+projekt

- Lab1: Grundläggande SQL
- Lab2: Databasdesign och EER modellerings
- Lab3: Avancerad SQL
- Lab4: Informationssäkerhet

Kursöversikt - LA+projekt

- 'Lab5': Projekt i bioinformatik
 - genomdatabas
 - proteindatabas
 - enzymdatabas
 - databas för biologiska reglersystem

Kursöversikt - LA+projekt

- Rapportatingsdeadline vid varje tentamenstillfälle
- behövs ett särskilt databaskonto
--> automatisk vid registrering på kursen
databaskontona tas bort efter 1 år
- anmälan till laborationer via kurshemsidan
senast 9 september

Examination

- tenta
- laborationsserie
- projekt

En kurs för TB

- Användning i senare kurser + arbete
- Unik och eftertraktad kompetens
 - Bio
 - Data
 - Förståelse av modellering + konsekvenser
(Hur modellera? Hur ställa frågor? Värför går det långsamt? Varför får man inget svar?...)

