

```
1 -----  
2     Database Systems (i.e. the worst core subject of the CS undergrad)  
3 -----  
4  
5 Data  
6 - Known facts stored and recorded  
7 - Complex objects such as text, numbers etc.  
8 - Raw information  
9 - e.g. a name or age of someone  
10  
11 Information  
12 - Data presented in context (can be summarised data)  
13 - Data that has been processed increasing the users knowledge.  
14 - e.g. "Matt de Bono started teaching Object Oriented Software Development in  
15     2017"  
16 // and really should've taught this class too  
17  
18 Data #is known and available  
19 Information #is processed and more useful  
20  
21 !What is a Database?:  
22 - A large, integrated, structured collection of data  
23 - Usually intended to model some real-world enterprise  
24  
25 !What is a Database System?:  
26 - A <Database Management System> is a software system designed to  
27     - Store  
28     - Manage  
29     - Facilitate access to databases  
30  
31  
32 Manipulate data with a query language (e.g. SQL)  
33  
34 #Types of Query Languages:  
35  
36 Data definition language (DDL):  
37 - Defines and sets up the database  
38  
39 Data manipulation language (DML):  
40 - Maintains and uses the database  
41  
42 Data control language (DCL):  
43 - To control access to the database  
44  
45 --- Problems you run into without Databases ---  
46  
47 # Program data dependence  
48 - If the <file structure changes>, so does the program  
49 - Programs "knows" too much about low-level data structure (access control)  
50  
51 # Redundant information  
52 - Wasteful, <inefficient>, loss of data integrity  
53 - Loss of metadata integrity  
54     - e.g. <same name different data>
```

```
55
56 # Limited Data Sharing
57 - Data is <tied to application> --> <hard/slow to create adhoc reports>
58
59 # Lengthly development times:
60 - Too much dependence on the application means the application has to do a lot
61   of low level data management, figure out file format each time
62
63 # Excessive program maintenance
64 - Up to 80% of development time in traditional file based organisations is for
65   maintenance
66
67 // and also no way to use the greatest programming language known to mankind:
68 // spicy query language
69
70 --- Good things about Databases ---
71
72 # Data Independence
73 - Separation of data and program, application logic
74 - Central data repository, central management
75
76 # Minimal Data Redundancy
77 - Redundancy can be controlled via <normalisation>
78
79 # Normalisation
80 - Breaking down a relationship between two relations into a set of smaller
81   relations
82
83 # Improved data sharing
84 - Data is shared
85 - External users can access
86 - Multiple views of data
87
88 # Reduced program maintenance
89 - Data structure can change without application changing
90
91 # Increased productivity of application development
92 - Data <already collected> and <structures already known>
93 - DBMS provides many tools to help access and manipulate data
94
95 # Enforcements of standards
96 - Centralised data management
97 - Documented policy for data management
98 - Data definition and dictionary (metadata)
99
100 # Improved data quality
101 - Constraints <built into the database>
102
103 // and a dead end career in a 9-5 enterprise job that you hate
104
105 So now with databases we can manage data in a structured way
106
107 Using the <relational model>
108 - <Rows> and <Columns> form <relations>
109 - <Keys> and <Foreign Keys> to link <relations>
110 // Basically just tables and headings like in spreadsheets
```

```
111
112 !Data Model:
113 - Collection of concepts for describing data
114
115 !Schema:
116 - Description of a particular collection of data, using a given data model
117
118 !Relational model of data:
119 - Most widely used model
120 - Relations have schemas which describes columns and fields
121
122 !Views:
123 - Describe how users see the data
124
125 !Conceptual Schema:
126 - A conceptual schema defines the logical structure
127
128 !Physical Schema:
129 - Describes the files and indexes used.
130
131 -----
132             Database Application Development Lifecycle
133 -----
134
135 Literally that thing from High School
136 For this subject, we will only focus on <design>, not any of the maintenance /
137 reevaluation stuff
138
139 # Database Design
140 - Conceptual
141 - Logical
142 - Physical
143
144 !Conceptual Design:
145 - Construction of a model of the data used in the database
146 - Independent of all physical considerations
147 - ER/EER diagrams
148
149 !Logical Design:
150 - Construction of the model based on the conceptual design
151 - Based on conceptual models
152 - Independent of a specific database and other physical considerations
153 - Decisions like what type of database should be used.
154 - Finding constraints
155     - What datatypes to use?
156
157 !Physical Design:
158 - Actually physically making the database
159 - Description of the implementation of the logical design for a specific DBMS
160 - Describes:
161     - Basic relations
162     - File organisations
163     - indexes
164     - Integrity constraints
165     - Security measures
166
```

```
171 # Guide to Conceptual Design
172 - What are the <entities> and <relationships> in the enterprise?
173 - what <information> about these entities and relationships should we store in
174   the database?
175 - What are the <integrity constraints> that hold?
176 - A database "schema" in the ER model can be represented as a ER diagram
177 - Can map an ER diagram into a relational schema
178
179 # ER Model Basics
180 !Entity:
181 - Real world object distinguishable from other objects
182 - An entity is described using a set of <attributes>
183
184 !Entity Set:
185 - A collection of similar entities, e.g. <all employees>
186 - All entities in an entity set have the same set of attributes
187 - Each entity set has a <key> (underlined)
188 - Each attribute has a <domain>
189
190
191 !Relationship
192 - Association among two or more entities (e.g. Fred works in Pharmacy dept)
193
194 !Relationship Set
195 - Collection of similar relationships
196
197 Same entity set can participate in different relationship sets, or in different
198 roles in the same set.
199
200 # Design Choices:
201 - Should a concept be modeled as an <entity or an attribute>?
202 - Should a concept be modeled as an <entity or a relationship>?
203 - Identifying relationships: <Binary or ternary>?
204
205 # Constraints in the ER Mode
206 - A lot of data semantics can and should be captured
207
208 --- Summary of Conceptual Design ---
209 - Yields a high level description of data to be stored
210
211 # ER model popular for conceptual design
212 - Constructs are expressive, and close to the way people think about their
213   applications
214 - Originally proposed by Peter Chen, 1976
215
216 # Basic Constructs:
217 - <Entities>, <Relationships> and <Attributes>
218 - <Weak Entities>
219
220 -----
221 Key Constraints
222 -----
```

223  
224 **!Participation Constraint:**  
225 - The participation of Departments in Manages is said to be **<total>** (vs partial)  
226 - When a attribute in one relation **<must appear in another relation's rows>**.  
227 e.g. A department relation with the attribute manager id also must appear  
228 in the manager relation.  
229 - Basically means "at least one"  
230

231 **!Weak Entity:**  
232 - Can be identified uniquely **<only by considering the primary key of another>**  
entity  
233 - Owner entity set and weak entity set must participate in one-to-many  
234 relationship set  
235 - Weak entity set must have total participation in this identifying relationship  
236  
237 Weak entities only have a "partial key" (dashed underline)  
238  
239

240 **Entity vs. Attribute:**  
241 - Should we model an address of an Employee to be an attribute or an entity?  
242 - **<Depends upon how we want to use address information>** and the semantics of the  
243 data  
244 - If we have **<several addresses per employee>**, address must be an entity  
245 - If the **<structure is important>** (e.g. city, street, etc.), address should  
246 be modeled as an entity  
247

---

248 -----  
249

250 ER design is **<subjective>**, there are many ways to model a given scenario  
251

252 **<Relational Database>** - A set of **<relations>**

253

254 **Relation:**  
255 - Made up of two parts, the **<Schema>** and the **<Instance>**

256

257 **Cardinality:**  
258 - Rows

259

260 **Degree / Arity:**  
261 - Fields

262

263 Can think of a relation as a set of rows or tuples  
264 - i.e. all rows are distinct, no order among rows

265

266 **Basics of SQL:**

267 [sql]

268 CREATE TABLE **<name>** (**<field>** **<domain>**, ...)

269

270 **INSERT INTO <name>** (**<field names>**)  
271     **VALUES** (**<field values>**)

272

273 **DELETE FROM <name>**  
274     **WHERE** **<condition>**

275

276 **UPDATE <name>**  
277     **SET** **<field name>** = **<value>**  
278     **WHERE** **<condition>**

```
279
280 SELECT <fields>
281     FROM <name>
282     WHERE <condition>
283 [/sql]
284
285 !Primary Key:
286 - A set of fields is a <super key> if
287     - No two distinct tuples can have the same values in all key fields
288     - Super key is like when you use <multiple attributes> to make the Keys
289
290 - A set of fields is a key for a relation if:
291     - It is a super key
292     - No subset of the fields is a superkey
293
294 What if there is >1 key for a relation?
295 - One of the keys chosen to be the <primary key>, other keys are called
296     <candidate keys>
297
298 It's *redundant* to have many <candidate keys unique> and have <one primary key>
299 Better just to have a <composite primary key>
300
301 !Foreign Key:
302 - Set of fields in one relation that is <used to refer to a tuple> in another
303     relation
304 - If all foreign key constraints are enforced, <referential integrity> is
305     achieved
306
307 Translating Conceptual to Logical:
308 Add constraints and keys to ER
309 - Underline           = PK      (Primary Key)
310 - Italic and Underline = FK      (Foreign Key)
311 - Underline and Bold   = PFK     (Primary Foreign Key)
312
313 From Logical to Physical Design:
314 Just make it in SQL lol
315
316 Logical to Physical Example:
317 Contracts (supplier_id, part_id, department_id) <- Imagine these are bold
318                                         and underlined (PK)
319 ---
320 translates to
321 ---
322 [sql]
323 CREATE TABLE Contracts (
324     supplier_id INTEGER,
325     part_id INTEGER,
326     department_id INTEGER,
327     PRIMARY KEY (supplier_id, part_id, department_id),
328     FOREIGN KEY (supplier_id) REFERENCES Suppliers,
329     FOREIGN KEY (part_id) REFERENCES Parts,
330     FOREIGN KEY (department_id) REFERENCES Departments)
331 [end]
332
333 Summary of the Relational Model:
334 - Tabular representation of data
```

```
335 - Simple
336 - Integrity constraints
337
338 // Skip lecture 5 because there's no way they could test us on how to workbench
339 // Note: I am not responsible for any damages to H1 dreams caused by this
340 // statement
341
342 // what is this subject structure
343
344 # Unary and Ternary Relationships
345
346 One-to-One and One-to-Many require <one foreign key> in the corresponding
347 relation
348
349 Many-to-Many you should generate an <associative entity> and <put two foreign>
350 <keys in the associative entity> (so you end up with two one-to-many relations)
351
352 # EER (Enhanced Entity Relation) Design
353
354 Basically most of database design depends on <how much you want to abstract> the
355 data, this whole thing is subjective and dependent on context.
356 // holy wow this subject is a meme
357
358 e.g.
359 For a database for a company that sells vehicles and requires a bunch of
360 different attributes, you could:
361 - separate it into the subtypes of vehicle (kinda like in OOP)
362 - or have it all at once in one relation (because you're a mad lad)
363
364 You'll need to explain why you made your choice!
365
366 For the first solution you could say <better query speeds> since less has to be
367 loaded into memory per relation.
368 Or for the second solution you could say <better space storage>.
369
370 Usually speed vs size or some weird query that you have to fulfil
371
372 Supertypes and Subtypes:
373 - Work the <same way as they do in OOP>
374 - <Supertype acts> kindof <like an super class>
375 - Subtypes inherit attributes (the same way subtypes extend supertype)
376 - Can be <disjointed> (only one subtype per instance) or <overlapping>
377
378 Crow's feet notation:
379 box with (d) - OR (d) =
380 Disjointed with one dash means that the instance can have <no subtype> if
381 necessary
382
383 Chen's:
384 double line to circle containing D or 0
385
386 Note that links go from one subtype to a super type means that there is <only a>
387 <link between that specific subtype> (it is possible to have a subtype that is
388 not related to a supertype).
389
390 // time for RELATIONAL ALGEBRA
```

## Relational Algebra

```
391
392
393
394
395 Query Language /= Programming Language
396
397 Relational Algebra:
398 - More operational
399 - Useful for representing execution plans
400
401 Relational Calculus:
402 - Describes what you want instead of how to compute it
403
404
405 !Selection ( $\sigma$ ):! // Unicode isn't playing nice so I'll use o
406 - Selects a subset of <rows> from relation
407
408 °rating>8(s2) // gets everything in s2 with rating above 8
409
410
411 !Projection ( $\pi$ ):
412 - Retains only wanted <columns> from relation
413
414 You use these by raising them to supertype ---  
for example, if I wanted to extract "sname" and "rating"  
from a relation called <s2>:
415
416 // for some reason I can't do superscript  $\pi$  and  $\sigma$ . :(
417
418 pisname,rating(s2) // gets sname and rating from s2
419
420 or if I wanted to get age:
421
422 piage(s2) // gets age from s2
423
424
425 # You can combine operations using brackets:
426
427 pisname,rating(°rating>8(s2))
428
429 If you can, try to be more specific with which table you are projecting from  
because then it is more efficient (less indexes to search)
430
431
432
433
434 !Set-difference (-):
435 - Tuples in r1, but not in r2 (as in the difference between the two relations)
436
437 To get the intersection: r1 - r2
438
439 !Union ( $\cup$ ):
440 - Tuples in both r1 or r2
441
442 To get the union: r1  $\cup$  r2
443
444 !Intersection ( $\cap$ ):
445 - Tuples that exist in both r1 and r2
446
```

```

447 To get the intersect: r1 n r2
448 -----
449 !Cross-Product (×):
450 - Allows you to combine two relations
451
452 S1 × R1: each row of S1 <paired with> each row of R1
453 (This is basically getting all permutations of every row)
454 [sql]
455
456 SELECT * FROM table1, table2;
457
458 [end]
459
460
461 ***Joining***
462
463 !Natural Join (⋈):
464 - Most common join, if one field (key) appears in both relations, then protect
465   all common attributes and copy each of the unique ones:
466
467 Example:
468 TableA                               TableB
469 Column1    Column2                  Column1    Column3
470 1          2                      1          3
471
472 TableA⋈TableB =
473
474 [sql]
475 SELECT * FROM TableA NATURAL JOIN TableB
476 [end]
477
478 column1  column2  column3
479 1          2          3
480
481 !Inner Join (or condition join):
482 - Similar to natural join but puts the other row next to the row with the same
483   specified condition
484
485 Example:
486 TableA                               TableB
487 Column1    Column2                  Column1    Column3
488 1          2                      1          3
489
490 TableA ⋈ e TableB = σe(R×S) // Where condition e is true
491
492 [sql]
493 SELECT * FROM TableA INNER JOIN TableB USING (Column1)
494 SELECT * FROM TableA INNER JOIN TableB ON TableA.Column1 = TableB.Column1
495 [end]
496
497 a.column1  a.column2  b.column1  b.column3
498 1          2          1          3
499
500 !Outer Join:
501 - Also works like inner join, but in cases where values do not exist for the
502   same key in one of the tables,

```

```
503
504 You can also rename things by using an arrow:
505 TableA(1 -> a1)
506 or by using px(E) <- where the Expression E is saved with the name x
507
508 -----
509                               What is SQL?
510 -----
511
512 SQL or SEQUEL is a query language used in relational database
513
514 The four basic commands are Create, Select, Update and Delete
515
516 DDL:
517 To set up/define database
518 - CREATE
519 - ALTER
520 - DROP
521
522 DML:
523 To maintain and use database
524 - SELECT
525 - INSERT
526 - DELETE
527 - UPDATE
528
529 DCL:
530 To control access to the database
531 - GRANT
532 - REVOKE
533
534 # Creating a table
535 If I wanted to create the relation Customer with the following attributes,
536 CustomerID, CustFirstName, CustMiddleName, CustLastName, BusinessName, CustType
537
538 You would do the following:
539 [sql]
540
541 CREATE TABLE Customer (
542     CustomerID      smallint          auto_increment,
543     CustFirstName   varchar(100),
544     CustMiddleName  varchar(100),
545     CustLastName    varchar(100)        NOT NULL,
546     BusinessName   varchar(200),
547     CustType        enum('Personal', 'Company') NOT NULL,
548     PRIMARY KEY (CustomerID)
549 );
550
551 [end]
552
553 This specifies that CustomerLastName and CustType cannot be null - they are
554 required in order to enter that row into the database.
555
556 CustomerID is simply a automatically incremented integer. It is also the primary
557 key
558
```

```
559 [sql]
560
561 CREATE TABLE Account (
562     AccountID      smallint      auto_increment,
563     AccountName    varchar(100)   NOT NULL,
564     OutstandingBalance DECIMAL(10,2) NOT NULL,
565     CustomerID    smallint      NOT NULL,
566     PRIMARY KEY (AccountID),
567     FOREIGN KEY (CustomerID) REFERENCES Customer(CustomerID)
568         ON DELETE RESTRICT
569         ON UPDATE CASCADE
570 );
571 [end]
572
573 In this case the <ON DELETE RESTRICT> and <ON UPDATE CASCADE> are used in
574 reference to any foreign keys.
575 - If there is an attempt made to delete in the parent table, it fails.
576 - If there is an update to the parent table then <the logical pointer in the>
577 <foreign key is also updated.>
578
579 # Aggregate Functions:
580 You can use these on a set of values in a column.
581 - AVG()
582 - MIN()
583 - MAX()
584 - COUNT()
585 - SUM()
586
587 All of these except count ignore null values.
588
589 e.g. SELECT COUNT(CustomerID) FROM Customer; //How many customers do we have
590
591 You can next things in SQL
592
593 Here's some comparison operators:
594 - IN / NOT IN
595 - ANY (true if any value returned meets condition)
596 - ALL (true of all values returned meet condition)
597
598 !Views:
599 A view is a relation that is <not in the conceptual/logical models> but is
600 <made available to the user> as a virtual relation is called a #view.
601
602 You can create views in the same way as you use relations:
603 [sql]
604
605 CREATE VIEW EmpPay AS
606 SELECT Employee.ID, Employee.Name, DateHired,
607     EmployeeType, HourlyRate AS Pay
608     FROM Employee INNER JOIN Hourly
609     ON Employee.ID = Hourly.ID
610 UNION
611 SELECT Employee.ID, Employee.Name, DateHired,
612     EmployeeType, AnnualSalary AS Pay
613     FROM Employee INNER JOIN Salaried
614     ON Employee.ID = Salaried.ID
```

```
615 UNION
616 SELECT Employee.ID, Employee.Name, DateHired,
617     EmployeeType, BillingRate AS Pay
618 FROM Employee INNER JOIN Consultant
619 ON Employee.ID = Consultant.ID;
620
621 [end]
622
623 You can actually insert values into views, but they won't be recorded in any
624 relation
625
626 // Go do some actual queries - reading these shitty notes doesn't really lend
627 // itself to the language since it's pretty interactive.
628 // W3 has some good tutorials.
629
630 // also this subject is by far the worst core one I've done so far
631 // We're about at lecture 9, 15 lectures to go
632
633 // the wild ride continues
634
635 -----
636                         Storage and Indexing
637 -----
638
639 !File:
640 A collection of <pages>, each containing a <collection of records>
641
642 A DBMS must support:
643 - Inserting/Deleting/Modifying records
644 - Reading a specified record
645 - Scanning all records
646
647 There are 3 main file organisations:
648
649 !Heap Files:
650 - <No order> among records
651 - Suitable when <all records> are <accessed at the same rate>, or you want all
652     to <retrieve all records>
653
654 This is the simplest data structure. As the file grows and shrinks, <disk pages>
655 are <allocated> and <de-allocated>.
656 #This data structure is the fastest for inserts.
657
658 In an implementation using lists, each page contains <data plus two pointers> to
659 <previous and next page>, just like <a linked list>.
660
661 The <header contains points to> the "head" (first full page) and "tail" (empty)
662
663 !Sorted Files:
664 - Records are <sorted by a certain condition>
665 - Best when some <records are required in a certain order>, or for retrieving a
666     <range of records>
667
668 Also works like heap files, i.e. a linked list, however the pages and records
669 <are ordered>.
670 This is faster for <searching for a specific record> because you can use
```

```

671 <binary search>
672 -----
673
674 !Comparison of Heap File vs Sorted File:
675
676 The <performance> of databases are normally <measured in Disk I/O's>
677
678 Where B == number of data pages:
679          Heap File      Sorted File
680 Scan All        B            B
681 1 Eq Search    0.5B (avg)    $\log_2 B$ 
682 Range Search    B             $(\log_2 B) + (\text{matching pages})$ 
683 Insert          2             $(\log_2 B) + 2*(B/2)$  // read to find pos, write
684 Delete          0.5B + 1     $(\log_2 B) + 2*(B/2)$ 
685
686 TLDR: Heap File Faster to Insert, Sorted File Faster to Search
687
688 -----
689
690 Sometimes we want to retrieve records by specifying values in one or more fields
691 e.g.
692 - Find all students in the "CIS department"
693 - Find all students with a gpa > 3
694
695 !Index File Organizations:
696 - "Special" Data structure
697
698 An <index> on a file speeds selections on the <search key fields> for the index
699 - A data structure built on top of data pages used for efficient search
700 - <Any subset> of the fields of a relation <can be the search key> for an index
701     on the relation
702 - <Search key> is <not the same as key>
703
704 The <index basically stores a bunch of pointers> to a <range of data entries> -
705 which themselves <link to the data record itself>
706
707 There are a few different types of Indexes:
708 - Clustered vs Unclustered
709 - Primary vs Secondary
710 - Single Key vs Composite
711 - Indexing Technique
712     - Tree-based, hash-based, other
713
714 Index Organisation Alternatives:
715 // this isn't explained so well
716 You can choose whether the index is linked to the record or the data itself
717 - whether you want variable record sizes, or just one bit of data per index
718 (this is the alternative 1, 2, 3 thing - alt 2 was used in assignment 3)
719
720 A good way to think of this is like keys to hash tables;
721
722 Alt 1:
723 - For an index, contain the actual data record
724
725 Alt 2:
726 - For an index, point to one record

```

727  
728 Alt 3:  
729 - For an index, point to a list of records  
730  
731 You can have <multiple indexes per file>  
732 - e.g. B+ Tree on age, Hash index on name, etc...  
733  
734 **!Clustered vs Unclustered:**  
735 - If the order of <data records> is the same as the order of <index data entries>  
736 then the index is called <clustered index>. Otherwise it is <unclustered>  
737 - A file can have a clustered index on at most <one search key>  
738 - Cost of retrieving data records through index varies greatly based on whether  
739 the index is clustered (cheaper for clustered)  
740  
741 **!Primary vs Secondary:**  
742 - Primary: index key includes the file's primary key  
743 - Secondary: has any other index  
744  
745 Primary index never contains duplicates (self explanatory since we're using  
746 primary key)  
747 Secondary index may contain duplicates.  
748  
749 **!Hash-based Index vs Tree-based index:**  
750 Hash-Based Index:  
751 - Good for <equality selections>  
752 - Work similar to an extendible hash table with buckets  
753 - Can't range search cause it would just spaghetti  
754  
755 Tree-Based Index:  
756 - Good for <range> selections  
757 - Hierarchical structure (tree) directs searches  
758 - Leaves contain data entries sorted by search key value  
759 - B+ tree: all root -> leaf paths have equal height (height)  
760 - So far we have seen tree-based indexes  
761  
762 Whats Examinable?  
763 - Describe the alternatives (1, 2, 3)  
764 - What is an index, when do we use them?  
765 - Index classification  
766  
767 -----  
768                           Query Processing (Selections & Projections)  
769 -----  
770  
771 Some queries are < expensive >  
772 So we need to actually think about what we're doing  
773  
774 ----- Selections -----  
775  
776 **!Simple Selections:**  
777 Of the form  
778 [sql]  
779 **SELECT \* FROM Reserves R WHERE R.name < "C%"**  
780 [end]  
781  
782 Size of Result:

```
783 size of R * reduction factor
784
785 #Reduction Factor or "Selectivity"
786 is how much you reduce the relation as defined in your condition.
787
788 //note: Reserves is of size 1000 records
789
790 With no index, unsorted:
791 - Must essentially scan the whole relation
792 - Cost is number of pages in R.
793 e.g. Reserves = 1000 I/O
794
795 With no index, sorted:
796 - Cost of binary search + number of pages containing results
797 - For reserves = 10 I/O + [selectivity * pages]
798
799 With an index:
800 - Use index to find qualifying data entries,
801 - Then retrieve the corresponding data records
802 Note: Hash index can only be used on equality
803
804 Cost Factors:
805 1. find qualifying data entries
806     (typically small, 2-4 I/O with B+ Tree, or 2.2 with hash index)
807 2. retrieve records
808
809 <Clustering helps a lot> with the <retrieving process>.
810
811
812 The approach to selecting records:
813 1. Find the cheapest path:
814 - an index or file scan with the fewest estimated page I/O
815
816 2. Retrieve the terms that match using this path
817
818 3. Apply the terms that don't match the index (if any)
819 - Other terms are used to discard some retrieved tuples
820
821 General Selection Conditions:
822 - First the selection is converted to <conjunctive normal form (CNF)>
823 e.g.
824 (day<8/9/94 AND rname='Paul') OR bid=5 OR sid=3
825 is converted to:
826 (day<8/9/94 OR bid=5 OR sid=3 ) AND (rname='Paul' OR bid=5 OR sid=3)
827
828 A B-Tree index matches a conjunction of terms that involve only attributes in a
829 <prefix of the search key>
830 - Index on <a, b, c> matches <a=5 AND b=3> but not <b=3>
831
832 Hash indexes must have all attributes in search key
833 // since you can't search for ranges in hash indexes
834
835 If we have 2 or more matching indexes:
836 1. get sets of rids of data records using each matching index
837 2. Then intersect these sets of rids
838 3. Retrieve the records and apply any remaining terms.
```

```
839
840 // I think this method is actually really similar to the set reduction used in
841 // declarative programming project 1? Minus the record retrieving bit.
842
843 Example: Consider name="Trump" AND is_president=True AND builds_wall=True
844 Assuming there is a B+ tree index on name and an index on builds_wall:
845 1. a) Retrieve rids of records satisfying "Trump" <using the first B+ index>
846     b) Retrieve rids of recs satisfying builds_wall=True <using the second index>
847 2. Intersect
848 3. Retrieve records and <check for is_president=True>
849
850 ----- Projection -----
851
852 !Projection with Sorting:
853
854 [sql]
855 SELECT DISTINCT R.sid, R.bid FROM Reserves R
856 [end]
857
858 How do we figure out duplicates?
859
860 Basic approach is to use <sorting>.
861
862 1. Scan R, extract only the needed attributes //just to save memory
863 2. Sort the resulting set
864 3. Remove adjacent duplicates
865
866 A cost estimation for this would be like this, assuming R has 1000 records:
867 Step 1: 1000 I/O
868 Step 2: 250 I/O          // Assuming Reserves has size ratio 0.25 = 250 pages
869                      // qsort is O(n log n) but we only calculate I/Os
870 Step 3: 250 I/O
871
872 This cost estimate assumes <everything fits into fast memory>.
873 But we can't always afford a few thousand terabytes of RAM :/
874
875 So, one way to get around this is with <external merge sort>
876
877 !Projection with Sorting Out of Memory (External Merge Sort):
878 This is a cool thing which allows you to use a sorting algorithm over a very
879 large amount of data.
880
881 Basically it works by sorting the data in parts.
882
883 These parts are then iterated parrallel to each other and merged, hence the name
884
885 Example of merge sort:
886
887 1. Split into smaller parts, then sort then (with like qsort or something)
888
889 2. Then look at all the parts parrallel to each other and merge
890
891 Where p1, p2 and out are buffers
892
893 p1           p2           out
894 | 1, 2, 5 |   | 1, 2, 3 |   | Empty |
```

```

895      ^
896
897 p1          p2          out
898 | 2, 5 |     | 1, 2, 3 |     | 1 |    // in this step, 1 already exists
899           ^                   // in output buffer, so discard 1
900
901 p1          p2          out
902 | 2, 5 |     | 2, 3 |     | 1 |
903           ^                   ^
904
905 p1          p2          out
906 | 5 |         | 2, 3 |     | 1, 2 |
907           ^                   ^
908
909 p1          p2          out
910 | 5 |         | 3 |     | 1, 2 |
911           ^                   ^
912
913 p1          p2          out
914 | 5 |         | ... |     | 1, 2, 3 | // once the buffer p2 runs out, load
915           ^                   // more data from the disk
916
917 p1          p2          out
918 | ... |       | ... |     | 1, 2, 3, 5 |
919           ^                   ^
920
921 Note for later:
922 As you can see, it takes 1 operation to sort the parts, and then another 1 to
923 merge into the out buffer. This makes up the "2" in the cost calculation later.
924
925 // I'm not going to explain it because the slides are spaghetti and this guy does
926 // a way better job:
927
928 // Watch them, they're only like 3 minutes each
929 ### https://www.youtube.com/watch?v=wTAVwbvjac
930 // I think this one is the one that's in the slides (about 1:30 onwards).
931
932 ### https://www.youtube.com/watch?v=1kr81S3HIR8
933 ### https://www.youtube.com/watch?v=5mKtLu60pdw
934
935 With external sort, the cost becomes this:
936 Cost with 20 buffer pages:
937 1. Scan R: 1000 I/O + 250 I/O // scanning/partitioning for external sort
938 2. Sort: 2 * 2 * 250 I/O
939     - The first "2" is for the sorting/merging, which is <done in 2 phases>
940     250 pages / 20 buffer -> 13 runs in order to fully sort/merge
941     <Remember we only count I/O, so the 2 represents one sort op and one >
942     <merge op>
943     - The second "2" represents a read/write operation
944     - Finally do this to all the pages - 250
945 3. Remove Duplicates: 250 I/O
946
947 Total: 1000 + 250 + 2*2*250 + 250 = 2500 I/Os
948
949
950 !Projection with Hashing:

```

```
951
952 This is a lot more simple than above, simply hash data into buckets and remove
953 any adjacent duplicates from the buckets.
954
955 1. Scan R, extract the needed attributes
956 2. Hash data into buckets
957 3. Remove adjacent duplicates from buckets
958
959 Cost:
960 1000 + 250 + 250 = 1500 I/O
961
962 But again, we can't always fit it in memory,
963 so
964
965 !Projection Hashing out of Memory:
966 1. Partition data into B partitions with h1 hash function
967 - Hopefully now a partition will fit into memory.
968
969 2. Go through each partition and hash each page with the h2 hash function
970
971 3. If a oversized partition is encountered, recursively partition this with step
972    1.
973
974 4. Check for duplicates in each h2 bucket
975
976 Cost is weird:
977 1. Read: 1000 I/O
978 2. H1 Write: 250 I/O
979 3. Read H1 Partitions: 250 I/O
980
981 Total 1500 I/O
982 -----
983
984 Sort-based approach is standard, doesn't have the problem of <skewing> with hash
985 and the <result is sorted> which is a nice bonus.
986
987 If there is <enough memory>, both have <same I/O cost>, as calculated by
988
989 M + 2T
990
991 Where:
992 M = pages in R
993 T = pages in R with unneeded attributes removed.
994
995 -----
996                               Query Processing (Joins)
997 -----
998
999 !Joins:
1000
1001 - Really common operation
1002 - Can be super <expensive> like cross products.
1003
1004 Here's some implementation techniques:
1005 - Nested-loops join
1006 - Index-nested loops join
```

```

1007 - Sort-merge join
1008 - Hash join
1009
1010 For the following examples, we have two Relations S and R.
1011 S - tuples are 50 bytes long, 80 tuples per page, 500 pages
1012     N = 500, ps = 80
1013 R - tuples are 40 bytes long, 100 tuples per page, 1000 pages
1014     M = 1000, pr = 100
1015
1016 R will be the outer relation and S will be the inner relation.
1017
1018 ----- Nested Loop Joins -----
1019
1020 !Simple Nested Loops Join:
1021 - literally in the name, <for each item, iterate through the relation and check>
1022
1023 Pseudocode:
1024 for i in R: // pr*M = 100*1000 Because we load every tuple
1025     for j in S: // N = 500 Here we only load the page
1026         if i==j: // once into fast memory
1027             result.append(i,j) // M = 1000 Output
1028
1029 Cost in I/Os:
1030 (pr * M) * N + M = 100*1000*500 + 1000
1031 ^ note that this is <the outer relation>, so it could be ps if S was outer.
1032
1033 !Page-Oriented Nested Loops Join:
1034 - Basically do simple nested loops join, but on <every single page in both R & S>
1035
1036 Pseudocode:
1037 for r_page in R: // M = 1000
1038     for s_page in S: // N = 500
1039         for i in r_page: // These don't count because we loaded the pages
1040             for j in s_page: // into fast memory
1041                 if i==j:
1042                     results.append(i,j) // M = 1000
1043
1044 This is <cheaper than doing a simple nested loop join>, as only the corresponding
1045 page needs to be loaded instead of the whole relation.
1046
1047 Cost:
1048 M*N + M = 1000*500 + 1000
1049
1050 !Index Nested Loops Join:
1051 // Note: Apparently this won't be examined because they messed it up in L14
1052
1053 - This is basically a simple nested loop join, HOWEVER it <utilizes an index> if
1054     the relation has one.
1055
1056 Pseudocode:
1057 for i in R: // pr*M = 100*1000
1058     j = index_search(i, S) // This can be 1.2 for hash index, or 2-4 for B+ tree
1059     result.append(i,j) // M = 1000
1060
1061 Cost:
1062 M + (pr*M)*(whichever index you're using + 1) = 1000 + 100*1000*2.2

```

```
1063  
1064 #In this case, the 2.2 comes from  
1065  
1066 1.2 I/O to find the record i in S, and 1 I/O to get the matching j record.  
1067 // using a hash index  
1068  
1069  
1070 !Block Nested Loops Join:  
1071 - This is really similar to page-oriented join, however multiple pages of R are  
1072     held, which is scanned against S and then output in sections called "blocks"  
1073  
1074 ----- Sort-Merge and Hash Joins -----  
1075  
1076 !Sort-Merge Join:  
1077 - Sort R and S.  
1078 - Then "merge" them, similar to <the method as described on line 890>, except  
1079     <only checking for equality instead of discarding equal values>.  
1080  
1081 Pseudocode:  
1082 sort(R)  
1083 sort(S)  
1084 merge R and S where r==s  
1085  
1086 Cost:  
1087 Sort R + Sort S + (M + N)  
1088  
1089 - <Less> sensitive to <data skew>  
1090 - <Result is sorted>  
1091 - <Goes faster if> one or both inputs <already sorted>  
1092  
1093 !Hash Join:  
1094 - Partition the relation the same way as <described on line 950>, except instead  
1095     <only checking for equality instead of discarding equal values>.  
1096  
1097 Cost:  
1098 Partitioning = 2(M+N)  
1099 Matching = M+N  
1100  
1101 Total = 2(M+N) + M+N  
1102  
1103 - <Better> if relation <sizes differ greatly>  
1104 - Highly <parallelizable>  
1105  
1106 ----- General Joins -----  
1107  
1108 So far we've done NL for equality conditions.  
1109  
1110 Inequality conditions <must use a B+ tree index>  
1111 # Hash Join, Sort Merge Join cannot be used for inequality!!!!  
1112  
1113 <Block Nested Loop is best for inequality conditions>  
1114  
1115 !Special Joins:  
1116 - Intersection  
1117 - Cross products  
1118 - Union (Distinct)
```

```
1119
1120 For Union:
1121 Sort-merge again
1122 - Sort both relations, scan sorted relations and merge
1123
1124 OR
1125
1126 Hash based
1127 - Repeat hash based process as described on line 950, but discard duplicates and
1128 add tuples to output.
1129
1130 !Aggregate Operations:
1131 - An aggregate option is an action that results in a single value e.g.
1132 AVG, SUM, MIN etc.
1133
1134 Without Grouping (like by a category):
1135 - Aggregate functions <require scanning the relation>
1136 - Can be done with an <index only scan>
1137
1138 With Grouping:
1139 - <Sort on group by attributes>, then <scan relation> and <compute aggregate>
1140 for each group.
1141
1142 -----
1143                         Query Optimization
1144 -----
1145
1146 <Many methods of executing> a query, all giving the same answer.
1147
1148 !Query Plan:
1149 This is a <tree> with <relational algebra operators for the nodes>, showing
1150 possible <choices of algorithms>.
1151
1152 By convention, <outer is on the left>.
1153
1154
1155 # How do I get gud and optimize f a s t?:
1156 1. Break query into "blocks"
1157     - Query Block = Unit of Optimization
1158     - Nested Blocks are treated as calls to a subroutine, made once per outer
1159         tuple
1160
1161 2. Convert each block <into relational algebra>
1162
1163 3. Consider all the <query plans> for <each block>
1164     - <Every query> has a <core SPJ select-project-join> expression
1165     - Group By, Having, Aggregation and Order By are <all done after SPJ>
1166     - <Focus on optimizing SPJ core>
1167
1168 4. Plan with the <lowest estimated cost is selected>
1169
1170 !Equivalence:
1171
1172 Selection <can be done condition by condition>; that is, if I ask for a name and
1173 age, you can <return age and then name in any order>
1174
```

```
1175 Projection, however, <must be done with the specified conditions all at once>
1176
1177 !De-Correlation:
1178 - Convert correlated subquery into uncorrelated subquery
1179
1180 Basically translating subqueries so that if they are in a nested for loop, they
1181 only need to be executed once:
1182
1183 [sql]
1184 SELECT S.sid
1185 FROM Sailors S
1186 WHERE EXISTS
1187   (SELECT *
1188     FROM Reserves R
1189     WHERE R.bid=103
1190     AND R.sid=S.sid) /* This is redundant */
1191
1192 /* vs */
1193
1194 SELECT S.sid
1195 FROM Sailors S
1196 WHERE S.sid IN
1197   (SELECT R.sid
1198     FROM Reserves R
1199     WHERE R.bid=103)
1200 [end]
1201
1202 !Flattening:
1203 Basically using a join algorithm + optimizer instead of two select-projects
1204 e.g.
1205
1206 [sql]
1207 SELECT S.sid
1208 FROM Sailors S
1209 WHERE S.sid IN
1210   (SELECT R.sid
1211     FROM Reserves R
1212     WHERE R.bid=103)
1213
1214 /* vs */
1215
1216 SELECT S.sid
1217 FROM Sailors S, Reserves R
1218 WHERE S.sid=R.sid
1219 AND R.bid=103
1220 [end]
1221
1222 !Cost Estimation:
1223 For each plan, estimate cost:
1224 - Estimate the cost of each operation in plan tree
1225 - Estimate the size of result for each operation
1226
1227 Performance Estimation:
1228 # I/O + factor * CPU instructions
1229
1230 Size Estimation:
```

```
1231 # Result = Max Number of Tuples * product of all RF's
1232
1233 RF = Reduction Factor, where
1234
1235 for equality, RF = 1/Keys(I)
1236 for range,    RF = (High(I)-value)/(High(I)-Low(I))
1237
1238 Assume RF = 1/10 otherwise
1239
1240 Estimated Size for Joins:
1241 Depending on which relation has more tuples,
1242
1243 est_size = NTuples(R) * NTuples(S) / MAX(NKeys(S), NKeys(R))
1244
1245 // Remember,
1246
1247 // What is query optimization/describe steps?
1248 // Equivalence classes
1249 // Result size/cost estimation
1250
1251 // is examinable and is probable a big part of it.
1252
1253 // try not to die
1254
1255 !Enumeration of Alternate Plans:
1256
1257 Single-relation plans
1258
1259 Multiple-relation plans
1260
1261 Single Relation Plans:
1262
1263 -----
1264                               Normalisation
1265 -----
1266
1267 All about the process of <abstraction> and <removing undesired redundancy>.
1268
1269 Not normalising tables will cause <anomalies>.
1270
1271 <Anomalies are not good!>
1272
1273 !Anomalies:
1274 # Insertion Anomaly
1275     e.g. A new course <cannot be added> until at least <one student has enrolled>
1276         - Which comes first? Student or course?
1277 # Deletion Anomaly
1278     e.g. If a <course only has one student>, and that <student withdraws>, we
1279         <lose all information regarding the course he took>!
1280 # Update Anomaly
1281     e.g. If the <student appears multiple times> and he changes year level, then
1282         we have to change <all the records containing that student>, otherwise
1283         we will <have an update anomaly>, where the data regarding an entity is
1284         <not certain>.
1285
1286 Basically just look for <repeating data> or <data that represents a separate entity>
```

1287  
1288 **!Functional Dependency:**  
1289 - A **functional dependency** concerns **values of attributes in a relation**.  
1290 - A set of attributes **X determines** another set of attributes **Y** if each value  
1291 of **X is associated with only one value of Y**.  
1292  
1293 **Written as:**  $X \rightarrow Y$  // Similar to  $y = f(x)$   
1294 EmployeeID → Name  
1295 EmployeeID → Salary  
1296  
1297 X is a **determinant** (an attribute on the left hand side of the arrow)  
1298  
1299 As stated before, all attributes are **either part of the primary key** **or not**.  
1300  
1301 A functional dependency can be **partial** if **one or more non-key attributes**  
1302 have **a dependency on part of the primary key**.  
1303 e.g.  
1304 If the dependency is  $X, Y \rightarrow Z$   
1305 then the partial dependence is  $Y \rightarrow Z$  or  $X \rightarrow Z$   
1306  
1307 A functional dependency is **transitive** if **it is a dependency between 2 or more**  
1308 **non-key attributes**.  
1309  
1310 ----- How to Normalise Real Good -----  
1311  
1312 **!First Normal Form (1NF): Remove Repeating Groups**  
1313  
1314 **Repeating groups of attributes shouldn't** be represented **in a flat table.**  
1315 Removing cells with multiple values will solve this!  
1316 e.g.  
1317 // Let's start with something flat. I'll use <> for underline  
1318 Order-Item(<OrderID>, CustomerID, <ItemID>, Desc, Qty)  
1319 // Notice how item has it's own attributes?  
1320 Order-Item(<OrderID>, CustomerID, (<ItemID>, Desc, Qty))  
1321 // Split into two  
1322 Order-Item(<OrderID, ItemID>, Desc, Qty)  
1323 Order(<OrderID>, CustomerID)  
1324 // Yay we're kinda normal now  
1325  
1326  
1327  
1328 **!Second Normal Form (2NF): Remove partial dependencies**  
1329  
1330 A **non-key attribute** **cannot be identified** by **part of a composite key**.  
1331  
1332 e.g. //cont.  
1333 Order-Item(<OrderID, ItemID>, Desc, Qty)  
1334 // Notice how Desc is identified by ItemID, which is part of the composite key.  
1335 Order-Item(<OrderID, ItemID>, Qty)  
1336 Item(<ItemID>, Desc)  
1337  
1338 Note that after we've done this, **we've solved most of the anomalies**.  
1339 1. Updating an **Item** will change the Desc wherever the Item is referenced  
1340 2. Data for an **Item** is **not lost when the last order for that item is deleted**  
1341 3. You can **add items without having to make an order for them**  
1342

1343  
1344  
1345 **!Third Normal Form (3NF): Remove transitive dependencies**  
1346  
1347 A non-key attribute <cannot be identified> by another non-key attribute.  
1348 e.g.  
1349 Employee(<EmpID>, Ename, DeptID, Dname) // Look for foreign keys  
1350 // In this case, Dname is identified by DeptID, so we normalise it  
1351 Employee(<EmpID>, Ename, <DeptID>)  
1352 Department(<DeptID>, Dname)  
1353  
1354  
1355  
1356 // this isn't even my final form  
1357 **!Boyce-Codd Normal Form: Ensure determinates are candidate keys**  
1358 "<Every non-key attribute> must <provide a fact about the key>, the whole key,  
1359 and nothing but the key. (So help me Codd)"  
1360 // woah this lecturer is way more fun  
1361  
1362 e.g.  
1363 Allocation(<StudentID, Subject>, Teacher, GPA)  
1364 // GPA is only linked to StudentID, however if Teacher is made part of the key,  
1365 // then GPA provides information about the teacher and the student, since  
1366 // it assume a student can take a subject more than once  
1367 Allocation(<StudentID, Teacher>, Subject)  
1368 Assignment(<StudentID, Teacher, GPA>)  
1369 Department(<Teacher>, Subject)  
1370  
1371 -----  
1372  
1373 However, <there is a pay-off> to having all the relations being neat and tidy with  
1374 no anomalies:  
1375  
1376 <Denormalized tables> are <fast>, but <more work must be done> to keep the  
1377 database consistent.  
1378 Denormalization <can be used to improve performance of time-critical things>  
1379  
1380 Normalised relations however <contain a minimum amount of redundancy> and allows  
1381 users to <insert, modify and delete> freely <without errors>.  
1382  
1383 -----  
1384 Database Administration  
1385 -----  
1386  
1387 The "Database Administrator" role // Only unlocked in New Game+  
1388  
1389 4 responsibilities as part of the <DBA role>  
1390 Capacity Planning:  
1391 - Estimating <disk space> and <transaction load>  
1392 Performance Improvement:  
1393 - Common approaches  
1394 // phat gains  
1395 Security:  
1396 - Threats  
1397 - Web apps and SQL Injection // look this up 2 be a cool h4x0r  
1398 // also firearms training and CQC

1399 **Backup and Recovery:**  
1400 - Types of failures, responses to these, types of backups  
1401 - Other measures to **<protect data>**

1402

1403 Primarily concerned with "maintenance" / "ops" phases  
1404 - But really **<should be consulted>** during **<all phases of development>**

1405

1406 DBA can be made **<redundant by cloud-based DBMS>**

1407

1408 **!Capacity Planning:**  
1409 Process of **<predicting when future load levels will saturate the system>** and  
1410 determining the **<most cost-effective way>** of **<delaying system saturation>** as  
1411 much as possible.

1412

1413 **Need to consider:**  
1414 - Disk space requirements  
1415 - Transaction throughput  
1416 // how much alcohol you'll need when the servers go down

1417

1418 **Estimating Disk Space Requirements:**  
1419 - Treat database size as the sum of all table sizes  
1420 - where table size = number of row \* average row width

1421

1422 Need to know the size of values. This is dependent on the storage engine used.  
1423 Go look at the documentation!  
1424 To calculate the **<width of the row>**, simply **<add up all the datatype sizes>**

1425

1426 **!Estimate growth of tables:**  
1427 Pretty simple, just gather estimates during **<system analysis>**.  
1428 e.g. "The company sells 1000 products. There are 2 million customers who place  
1429 on average 5 orders per month. An average order is for 8 different  
1430 products".

1431

1432 will make:  
1433 Product -> 1000 rows  
1434 Customer -> 2,000,000 rows  
1435 Orders -> 10,000,000 rows  
1436 OrderItems -> 80,000,000 rows

1437

1438 Figure out storage per year according to the new data.  
1439 Note that "Event" tables grows a lot faster than "entity" tables, so we can put \  
1440 event tables on separate disks.

1441

1442 **!Estimate transaction load:**  
1443 - Consider each business transaction.  
1444 - How often are these run?  
1445 - For each transaction, what SQL statements will be run?

1446

1447 Multiply all of these and then you get the transaction load.

1448

1449 **Things to consider:**  
1450 - Differentiate **<peak vs average load>**.  
1451 e.g. a sale will result in **<many more transactions>**  
1452 - Acceptable **<response time>**  
1453 - We can **<reduce the response time>** by:  
1454 - Application/Database **<Design>**

1455 - <Tuning> DBMS/individual queries  
1456 - <not running everything on a potato>  
1457 - Availability  
1458 - 9-5 vs 24/7  
1459 - Who's the customer? e-business? Do they need this all the time?  
1460  
1461 What affects database performance?:  
1462 - Caching data in memory  
1463 - Placement of data files  
1464 - Database replication and server Clustering  
1465 - Use of fast storages such as SSDs  
1466 - Use of indexes to speed up searches/Joins  
1467 - You must decide <when to create indexes>  
1468 - Frequently queried columns  
1469 - Primary/Foreign Keys / Unique Columns  
1470 - Large tables  
1471 - Columns with a wide range of values  
1472 - Choice of data types  
1473 - Program logic  
1474 - Query execution plans  
1475 - Good code (no deadlocks)  
1476  
1477 **!Security:**  
1478  
1479 Threats to Databases:  
1480 Loss of Integrity:  
1481 - Keep data consistent  
1482 - Free of errors/anomalies  
1483  
1484 Loss of Availability:  
1485 - Must be available to authorized users for authorized purposes  
1486  
1487 Loss of confidentiality:  
1488 - Must be protected against unauthorized access  
1489  
1490 We can <protect databases> by:  
1491 - Access Control  
1492 - Encryption  
1493  
1494 **!Access Control:**  
1495 All <about restricting access> to data and <granting and revoking> privileges  
1496 This <is required by the security mechanism of a DBMS>  
1497  
1498 - Access control is handled by the admin creating user accounts  
1499 - Database keeps a usage log  
1500 - When <tampering is suspected, perform an audit>.  
1501 - Need to control <online and physical access> to the database  
1502  
1503 Types of Privileges:  
1504 - Account Level  
1505     DBA specifies the particular privileges that each user holds  
1506         i.e. this user <can do x> on the database  
1507 - Table Level  
1508     Controls a user's privileges to <access particular tables> or views  
1509 - Schema Level  
1510     Controls a user's privileges to <access a particular scheme> in the database

```
1511  
1512 Views are pretty important since you can <hide database structure and data>  
1513  
1514 You can also <encrypt> things to <protect sensitive data> whenever they are  
1515 transmitted over a network  
1516 - This <prevents interception by third party>  
1517  
1518 Encrypt data in the database  
1519 - Provides some protection in case of unauthorized access.  
1520  
1521 !SQL Injection:  
1522 // this is fun  
1523 A technique used to exploit web applications that use user input within  
1524 database queries.  
1525  
1526 Basically you can <run queries in fields like username and password>.  
1527  
1528 ### https://xkcd.com/327/  
1529  
1530 One example of a h4x0r injection  
1531  
1532 usr: ' or 1=1; --  
1533 pw: any text  
1534  
1535 will result in:  
1536 [sql]  
1537 SELECT * FROM User  
1538 WHERE username = '' or 1=1; -- ' and password = 'any text';  
1539 [end]  
1540 The ' closes the quotes, then there is an always true condition, and the rest  
1541 of the query is commented out.  
1542  
1543 !Backup and Recovery:  
1544 This is pretty self explanatory  
1545  
1546 backup ur shit boi  
1547  
1548 Types of Backup:  
1549 # Physical Backup  
1550 - Raw copies of Files and Directories  
1551 - Suitable for large databases that need fast recovery  
1552 - Database is preferably offline ("cold" backup) when backup occurs  
1553 - Backup = exact copies of files  
1554 - Backup <should include logs>  
1555  
1556 # Logical Backup  
1557 - Backup <completed through SQL Queries>  
1558 - Slower than physical  
1559 - <Output is larger than physical>  
1560 - <Doesn't include log or config files>  
1561 - Machine independent  
1562 - Server is available during the backup  
1563  
1564 You can backup Online - "HOT" or Offline - "COLD"  
1565  
1566 Online Backups are done <while the database is live>.
```

1567 Offline backups are done <while the database is offline>.  
1568 - Offline is preferable  
1569 - Simpler to perform  
1570 - Backups occur when the database is stopped  
1571  
1572 You can backup Fully or Incrementally.  
1573 Fully contains the <whole database> // livin in the database  
1574 Incremental <only contains changes> since the last backup  
1575  
1576 You should have a <backup strategy> which is a <combination of full/incremental>  
1577 backups  
1578  
1579 You should also have <offsite backups> just in case your building explodes or  
1580 something  
1581  
1582 Should also use <Server Replication>, <Server Cluster> and <RAID> // love RAID  
1583 // RAID 0 FOR FULL BACKUP POTENTIAL (not really)  
1584  
1585 **Replication:**  
1586 - One writer/Many readers  
1587 - Some protection against server failure  
1588 - Multiple copies of data  
1589 - Replicates accidental data deletion  
1590  
1591 **Clusters:**  
1592 - Automatic synchronous partition  
1593 - Protection against server failure  
1594 - Multiple copies of data  
1595  
1596 **RAID:**  
1597 - Software or Hardware RAID  
1598 - Depends on the RAID level // raid 0 will be fast but won't help backup  
1599  
1600 **!Categories of Failure:**  
1601  
1602 **Statement Failure:**  
1603 - Bad syntax  
1604  
1605 **User Process Failure:**  
1606 - The process doing the work fails  
1607  
1608 **Network Failure:**  
1609 - Network failure between the user and the database  
1610  
1611 **User Error:**  
1612 - User accidentally drops the rows, table, database  
1613 // me\_irl  
1614  
1615 **Memory Failure:**  
1616 - Memory fails, becomes corrupt  
1617  
1618 **Media Failure:**  
1619 - Disk failure, corruption deletion  
1620 // i.e. your computer catches fire  
1621  
1622 -----

1623  
1624  
1625  
1626 **!Transaction:**  
1627  
1628  
1629 <The system> should <ensure that> updates of a <partially executed transaction>  
1630      are <not reflected in the database>.  
1631  
1632 A good transaction requires #ACID  
1633  
1634 **!ACID:**  
1635 Stands for:  
1636 - Atomicity  
1637      - A transaction is a <single, indivisible, logical unit of work>.  
1638      - <All operations> in a transaction <must be completed>; if not, then the  
1639           transaction is aborted.  
1640  
1641 - Consistency  
1642      - <Data constraints> that hold <before> a transaction must also hold <after it>  
1643           e.g. PKs and FKs and implicit constraints (sums etc.)  
1644  
1645 - Isolation  
1646      - <Data used during execution> of a transaction <cannot> be used by a second  
1647           transaction <until the first one is completed>  
1648      - Solved by <running transactions serially> (one after another).  
1649  
1650 - Durability  
1651      - When the <transaction is complete> the changes made are <permanent> even  
1652           if the <system fails>  
1653  
1654 If a transaction has <multiple statements>, they are usually <preceded> by a  
1655 <START TRANSACTION> or <BEGIN> and <ended> by a <COMMIT>  
1656  
1657 [sql]  
1658 **START TRANSACTION;** -- (or, 'BEGIN')  
1659 SQL statement;  
1660 SQL statement;  
1661 SQL statement;  
1662 ...  
1663 **COMMIT;**  
1664 [end]  
1665  
1666 ROLLBACK is used to <undo everything> (in case the transaction doesn't finish)  
1667 Transactions are "all or none"  
1668  
1669 **!Concurrent Access:**  
1670 When <multiple users> execute DML (Database Modification Language) against a  
1671 <shared database>  
1672  
1673 Can result in <uncommitted data> if two queries are <executed at the same time>  
1674  
1675 e.g. If you change the bank account of someone <while they're withdrawing>, then  
1676       <commit the changes before the other person finishes>, then you <do not>  
1677       record changes of <both executions>.  
1678

1679 It'd be like **<item duping>** if one person held an item while another person drops  
1680 it from the inventory

1681

1682 **!Serialisation:**

1683 We can **<serialise>** a transaction by ensuring all executions are done **<one after>**  
1684 another

1685

1686 However this is **<real slow>** and **<not realistic>**

1687

1688 So we have **<concurrency control>**

1689

1690 **!Concurrency Control:**

1691 - Locking/Version Control

1692

1693 **!Locking:**

1694 - **<Guarantees exclusive use>** of a data item to a **<current transaction>**

1695 - **<Required to prevent>** another transaction from **<reading inconsistent data>**

1696

1697 The **<Lock Manager>** is required for assigning/policing the locks

1698

1699 **Database Level Lock:**

1700 - Entire database is locked

1701 - Good for batch processing but unsuitable for multi-user DBMSs

1702

1703 **Table Level Lock:**

1704 - Entire table is locked

1705 - Multiple users can use the database as long as they **<use separate tables>**

1706 - Can cause bottlenecks

1707 - Not suitable for highly multi-user DBMSs

1708

1709 **Page Level Lock:**

1710 - Page is locked

1711

1712 **Row Level Lock:**

1713 - **<Most Popular>**

1714 - Locks only the row of a table

1715 - Improves **<data availability>** but with **<high overhead>**

1716

1717 **Field Level Lock:**

1718 - Most flexible lock but requires **<extremely high level of overhead>**

1719

1720 Locking can result in a **<dead lock>**, where two transactions wait for each other  
1721 to unlock data.

1722

1723 Prevent this by locking potential records in advance - detect/end the transaction

1724

1725 **!Versioning:**

1726 - **<Timestamping>** all the transaction

1727 - Each **<transaction attempts an update>**

1728 - System will review and **<reject inconsistent updates>**

1729

1730 -----

1731 Data Warehousing

1732 -----

1733

1734 **!Data Warehouse:**

- A <single repository> of <organisational data>
  - Integrates data from multiple sources
  - Pretty much a <informational database>
- Transactional vs Informational:**
- a Transactional system's purpose is to <run the day to day business>, whereas the informational system's purpose is to <support decision making>
- A transactional database represents the <current state of business> whereas the informational database represents <historical data>.
- This is also reflected by the users - customers/employees will use transactional while <managers and analysts> use <informational>
- Dimension Modelling (Star Schema Design):**
- Consists of:
- A fact table
  - several dimensional tables
  - Hierarchies in the dimensions
- Essentially a <simple and restricted> type of ER model.
- Fact Table:**
- A fact table contains the actual business measures called facts
  - Also contain foreign keys for dimensions
  - Really really simple ER model
- 
- NoSQL
- 
- As much as the relational database we've learnt about is <flexible> and <fast>, relational dbs are <not good with big data> and <clustered/replicated servers>
- NoSQL Database:**
- Doesn't use relational model
  - Runs well on <distributed servers>
  - <Scaling out> rather than <scaling up>
  - <Open Source> usually
  - Built for <modern web> and <natural for cloud environment>
  - <Schema-less>
  - Not ACID compliant
- Built to <handle larger data volumes>
- Some examples of <NoSQL> are <JSON> and <XML>.
- Because NoSQL languages are <built around big data>, they just capture everything in a <data lake> and <worry about structure later>.
- Data Lake:**
- A large integrated repository for internal and external data that <doesn't follow a predefined schema>
  - Capture <everyting>

```
1791 Usually <structuring and organizing the data> takes place <during data analysis>
1792
1793 There are many different types of NoSQL:
1794 # Key-Value store
1795     e.g. JSON
1796 # Document Databases
1797     Examinable/Queriable key-value dbs
1798     e.g. MongoDB, CouchDB, Terrastore, OrientDB, RavenDB
1799 # Column Families
1800     Columns are stored <instead of rows>
1801     Kind of like automatic vertical partitioning
1802     e.g. DynamoDB, Cassandra
1803 # Graph Databases
1804     Node-and-arc network
1805     Difficult to program in relational DB
1806     Stores entities and their relationships
1807     e.g. Neo4j
1808
1809 !CAP Theorem:
1810
1811 PICK TWO:
1812 # Consistency:
1813     - Everyone always sees the same data
1814 # Partition Tolerance:
1815     - System stays up when network between nodes fail
1816 # Availability:
1817     - System stays up when nodes fail
1818
1819 Or alternatively, IF a partition occurs, <you must choose Consistency or Availability>
1820
1821 // godspeed guys
1822 // im pretty dead
1823 // o7
```