

# **Базы данных**

**Лекция 02 – Основные концепции. Терминология**

**Преподаватель: Поденок Леонид Петрович, 505а-5**

**+375 17 293 8039 (505а-5)**

**+375 17 320 7402 (ОИПИ НАНБ)**

**prep@lsi.bas-net.by**

**ftp://student:2ok\*uK2@Rwox@lsi.bas-net.by**

**Кафедра ЭВМ, 2022**

2022.09.15(17)

## Оглавление

|                            |   |
|----------------------------|---|
| Основная терминология..... | 3 |
|----------------------------|---|

# Основная терминология

Существует много вариантов терминов, используемых при описании данных.

Широко признанным авторитетом в области баз данных, не связанным ни с какой конкретной фирмой-производителем ЭВМ, является CODASYL (COference on DAta SYstems Languages — Ассоциация по языкам для систем обработки данных).

## Байт

Байт — наименьшая адресуемая группа битов (обычно состоит из 8 бит[ов]).

## Элемент данных

Элемент данных — наименьшая единица поименованных данных. Может состоять из любого количества битов или байтов. Часто элемент данных называют *полем*.

Элементом данных может быть некоторая величина, например, в счет-фактуре.

## Агрегат данных

Агрегат данных — поименованная совокупность элементов данных внутри *записи*, рассматриваемая как единое целое.

Например, агрегат данных **ДАТА** может состоять из элементов данных **МЕСЯЦ, ДЕНЬ, ГОД**.

Существуют два типа агрегатов данных – векторы и повторяющиеся группы.

**Вектор** — одномерная упорядоченная совокупность элементов данных (например, **ДАТА**).

**Повторяющаяся группа** — совокупность данных, которые встречаются несколько раз в экземпляре записи, например прием и выдача вкладов в записи счета сберкассы.

В повторяющуюся группу могут входить отдельные элементы данных, векторы, агрегаты данных или другие повторяющиеся группы.

## **Запись**

Запись — поименованная совокупность элементов или агрегатов данных.

Запись — агрегатный тип данных, инкапсулирующий набор значений различных типов без их сокрытия.

При чтении из базы данных может быть прочитана логическая запись как целиком, так и частично — в ряде случаев логической записью базы данных является структура данных, в которую входит несколько групп элементов данных (*сегментов*), не все из которых имеет смысл читать одновременно.

Верхний предел для числа возможных экземпляров записи конкретного типа зависит, в основном, от наличия необходимого оборудования для хранения и доступа.

Верхний предел для числа повторяющихся групп внутри записи обычно ограничен небольшим числом.

## **Сегмент**

Существует мнение, что нет необходимости различать агрегат данных и запись, поскольку и то и другое является совокупностью элементов данных. В терминологии, используемой фирмой IBM, а также в других источниках и агрегат данных, и запись называют сегментом.

Сегмент состоит из одного или нескольких элементов данных (обычно нескольких) и является основным квантом данных, передаваемым прикладной программе или получаемым от нее под управлением программного обеспечения, работающего с базой данных.

## **Физический уровень**

- уровень блоков;
- уровень файлов.

## Таблица

Записи одинаковой структуры обычно организуются в двумерные таблицы и представлены в ней строками. Элементы данных записей организованы в виде столбцов. Таблицы организованы в *базы данных*. Ниже БД, содержащая 2 таблицы, связанные через поле DEPT\_NO.

### DEPARTMENT – Отделы

| DEPT_NO | DEPT_NAME   | BUDGET        |
|---------|-------------|---------------|
| D1      | Marketing   | 5 000 000.00  |
| D2      | Development | 25 000 000.00 |
| D3      | Research    | 10 000 000.00 |

### EMPLOYEES – Служащие

| EMP_NO | EMP_NAME  | DEPT_NO | SALARY    |
|--------|-----------|---------|-----------|
| E1     | Ivanov    | D1      | 40 000.00 |
| E2     | Smirnov   | D2      | 45 000.00 |
| E3     | Kuznetsov | D1      | 41 000.00 |
| E4     | Popov     | D2      | 43 000.00 |
| E5     | Vasiljev  | D3      | 55 000.00 |
| E6     | Petrov    | D2      | 45 000.00 |

В разное время и в различных контекстах термин **запись** может означать экземпляр записи или тип записи, логическую запись или физическую запись, хранимую запись или виртуальную запись, а возможно, и еще что-нибудь. В связи с этим в формальной реляционной модели термин запись не используется — вместо него применяется термин *кортеж* (tuple). Также в реляционной модели не используется термин *поле*, вместо него используется термин *атрибут*.

## **Кортеж**

Термин *кортеж* в отношении к базам данных приблизительно соответствует понятию строки в таблице (так же, как термин *отношение* приблизительно соответствует понятию таблицы).

Соответственно, *кортеж* состоит из атрибутов.

*Кортежем* вообще называется набор взаимосвязанных величин.

*Кортеж*, содержащий две величины, называется *парой* (pair).

*Кортеж*, содержащий *N* величин, называется *N-кортежем*.

Таблица состоит из набора *кортежей*, каждый из которых содержит элементы данных одинакового типа. Она, таким образом, представляет собой двумерную матрицу элементов данных.

Элементы данных обычно обрабатываются по группам. В различных системах программного обеспечения эти группы называются по-разному (*сегмент*, *кортеж*). Общепринятым является термин *запись*.

## **Отношение**

*Отношение* — это формальное название таблицы.

Например, можно сказать, что база данных отделов и служащих, представленная выше, содержит два *отношения*.

В настоящее время в неформальном контексте термины *отношение* и *таблица* принято считать синонимами. На практике в подобном контексте термин *таблица* используется гораздо чаще, чем термин *отношение*.

## Файл

Файл — поименованная совокупность всех экземпляров логических записей заданного типа.

В простом файле в каждой логической записи содержится одинаковое число элементов данных (рисунок 1).

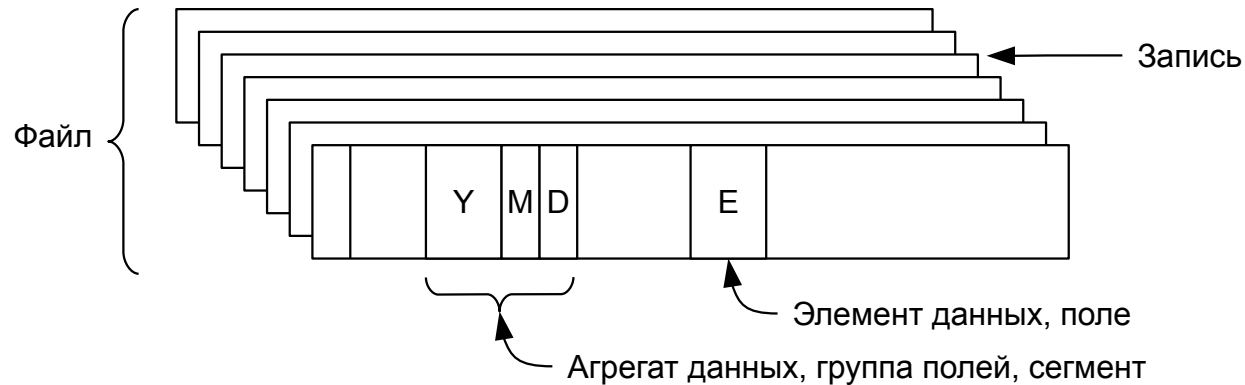


Рисунок 1 — Термины, используемые при описании данных с точки зрения прикладного программиста

В более сложном файле из-за наличия повторяющихся групп записи могут состоять из различного числа элементов данных.

## База данных

База данных — совокупность экземпляров различных типов записей и отношений между записями, агрегатами данных, элементами данных.

## **Пример – счет-фактура**

Некоторая величина в счет-фактуре является *элементом данных*.

*Агрегатом данных* может быть строка в счет-фактуре. Если она повторяется несколько раз, она будет *повторяющейся группой*.

Строка может рассматриваться, как отдельно адресуемая группа данных, в этом случае она будет *сегментом*.

Весь счет-фактура может рассматриваться, как *логическая запись*.

Вся совокупность записей счетов-фактур обычно хранится в одной или нескольких *таблицах*.

Физическая структура хранения логических записей на устройстве хранения зависит от СУБД.

## **Виртуальные данные**

Слово *виртуальный*, относящееся к техническим средствам или данным, указывает на то, что некоторый элемент в запросе представляется прикладному программисту существующим, тогда как фактически в представленном виде он отсутствует.

Программист может обращаться за виртуальными данными, которые им предполагаются существующими, но не существуют фактически, по крайней мере в данной форме. Каждый раз, когда программист обращается за этими данными, они генерируются определенным образом, что возможно обеспечивает более компактную и удобную форму их хранения.

## **Прозрачные средства и данные**

Виртуальное только представляется существующим, прозрачное представляется несуществующим, но на самом деле существует. Многие данные и механизмы, используемые при их хранении и передаче, могут быть скрыты от программиста (view).



## **База данных**

Базу данных можно определить как совокупность взаимосвязанных хранящихся вместе данных при наличии такой минимальной избыточности, которая допускает их использование оптимальным образом для одного или нескольких приложений.

Данные хранятся таким образом, чтобы они были независимы от программ, использующих эти данные.

Для добавления новых или модификации существующих данных, а также для поиска данных в базе данных применяются общие методы (способы).

Данные структурируются таким образом, чтобы была обеспечена возможность дальнейшего наращивания приложений.

Говорят, что система содержит совокупность баз данных, если эти базы данных структурно полностью самостоятельны.

База данных может разрабатываться для:

- пакетной обработки данных;
- обработки в реальном времени (жесткие ограничения на время обработки запроса);
- оперативной обработки (в этом случае обработка каждой транзакции завершается к определенному моменту времени, но при этом на время обработки не накладывается жестких ограничений, существующих в системах реального времени).

Во многих базах данных предусмотрена совокупность вышеуказанных методов обработки, причем все три вида обработки могут выполняться параллельно.

## Избыточность

База данных представляет возможность в значительной степени избавиться от избыточности.

Базу данных иногда определяют как *неизбыточную совокупность элементов данных*, однако в действительности избыточность в определенной степени допускают во многих базах данных с целью уменьшения времени доступа к данным и/или упрощения способов адресации.

Иногда некоторые записи дублируются для того, чтобы обеспечить возможность быстрого восстановления данных при их случайной потере.

Чтобы база данных была неизбыточной и удовлетворяла другим требованиям, приходится идти на компромисс. В этом случае говорят об *управляемой*, или *минимальной*, избыточности или о том, что хорошо разработанная база данных свободна от *излишней* избыточности.

Неуправляемая избыточность имеет несколько недостатков:

- 1) хранение нескольких копий данных приводит к дополнительным затратам;
- 2) что особенно серьезно, приходится выполнять многократные операции обновления для нескольких избыточных копий.

Избыточность поэтому обходится значительно дороже в тех случаях, когда при обработке файлов обновляется большое количество информации или, что еще хуже, часто вводятся новые элементы данных или уничтожаются старые.

- 3) вследствие того что различные копии данных могут соответствовать различным стадиям обновления, информация, выдаваемая системой, может быть противоречивой.

Если не использовать базы данных, то при обработке большого количества информации появится так много избыточных данных, что фактически станет невозможным сохранять их все на одном и том же уровне обновления.

## **Непрерывное расширение**

Одной из наиболее важных характеристик большинства баз данных является их постоянное изменение и расширение.

По мере добавления новых типов данных или при появлении новой бизнес-логики должна быть обеспечена возможность быстрого изменения структуры базы данных.

Реорганизация базы данных должна осуществляться по возможности без перезаписи прикладных программ и в целом вызывать минимальное количество преобразований.

Простота изменения базы данных может оказать большое влияние на развитие приложений для баз данных, используемых в управлении производством.

Требования к обработке обычно изменяются непредсказуемым образом. При этом если возникает необходимость модификации выбранных структур данных, то приходится соответственно перезаписывать и отлаживать прикладные программы. Чем большее количество прикладных программ имеется в наличии на установке, тем более дорогой становится эта процедура.

Поэтому одним из важных свойств базы данных является независимость данных и использующих их прикладных программ друг от друга в том смысле, что изменение одних не приводит к изменению других.

В действительности же полностью независимыми данные бывают так же редко, как и полностью избыточными.

Независимость данных определяется с различных точек зрения. Сведения, которыми должен располагать программист для доступа к данным, различны для различных баз данных.

Тем не менее независимость данных — это одна из основных причин использования систем управления базами данных.

## **Установление многосторонних связей**

В том случае, когда один набор элементов данных используется для многих приложений, между элементами этого набора устанавливается множество различных взаимосвязей, необходимых для соответствующих прикладных программ.

Организация базы данных в значительной степени зависит от реализации связей между элементами данных и записями.

В базе данных, используемой многими приложениями, должны быть установлены многочисленные промежуточные взаимосвязи между данными. В этом случае при хранении и использовании данных контролировать их правильность, обеспечивать их защиту и секретность труднее, чем при хранении данных в простых, несвязанных файлах. Что касается обеспечения секретности данных и восстановления их после сбоев, то этот вопрос является очень важным при конструировании баз данных.

В подавляющем количестве систем средства управления базами данных обеспечивают возможность пользователи использовать данные таким способом, который не был предусмотрен разработчиками системы. Это значит, что пользователи могут обращаться к БД с запросами, которые заранее в ней не предусматривались. Наличие такой возможности означает организацию данных в системе, при которой доступ к ним осуществляется по различным путям, причем одни и те же данные могут использоваться для ответа на различные запросы.

При организации данных в виде ориентированных на конкретное приложение схем, в которых доступ к данным производится одним и тем же определенным способом, описанная выше возможность отсутствует. Вся существенная информация об объектах должна сохраняться одновременно и полностью, а не только та ее часть, которая необходима для конкретного варианта использования.

## Идентификатор объекта

Программисту или администратору базы данных необходимо иметь возможность обращаться к записи или к кортежу, связанному с данным объектом.

Для этого необходимо иметь возможность идентифицировать запись или группу записей и располагать средствами их обнаружения в таблице.

С этой целью один из элементов данных обычно определяется в качестве идентификатора объекта.

Идентификатор объекта должен быть уникальным – никакой другой объект не может иметь то же значение данного элемента данных.

Идентификатором объекта **СЛУЖАЩИЙ** может быть **НОМЕР-СЛУЖАЩЕГО**.

Идентификатором объекта **СЧЕТ** может быть **НОМЕР-СЧЕТА**.

Иногда требуется более одного элемента данных для идентификации записи.

Например, для того чтобы идентифицировать банковскую операцию для записи **СЧЕТ**, необходимы идентификаторы **НОМЕР-СЧЕТА** и **НОМЕР-ОПЕРАЦИИ**.

Для идентификации записи о рейсах самолетов необходимы **НОМЕР-РЕЙСА** и **ДАТА**. Одного номера рейса недостаточно, так как вылет с одним и тем же номером рейса может происходить каждый день.

## Первичный ключ

Идентификатор объекта рассматривается как ключ записи или группы записей. Такой ключ называется *первичным ключом*.

Существуют также и вторичные ключи, которые не идентифицируют объект (единственным образом).

В тех случаях, когда в качестве ключа используются несколько элементов данных (полей, агрегатов), их обычно указывают как соединенные символом «плюс», например **НОМЕР-РЕЙСА + ДАТА**. Такое сочетание называется сцепленным ключом.

Иногда для формирования первичного ключа необходимы три, четыре и более элементов данных. Запись о кассовом сборе, например, в приведенной ниже группе идентифицируется комбинацией элементов **КИНОФИЛЬМ, КИНОТЕАТР, ДАТА**.

|                  |                  |             |                      |
|------------------|------------------|-------------|----------------------|
| <b>КИНОФИЛЬМ</b> | <b>КИНОТЕАТР</b> | <b>ДАТА</b> | <b>КАССОВЫЙ-СБОР</b> |
|------------------|------------------|-------------|----------------------|

Первичным ключом поэтому является **КИНОФИЛЬМ+КИНОТЕАТР+ДАТА**.

*Первичный ключ* — это такой элемент данных или такая совокупность элементов данных, которая единственным образом идентифицирует одну запись или группу записей.

Первичный ключ имеет большое значение, так как он используется для определения местоположения записи с помощью индексов или других методов адресации.

## Избыточные значения

Значения атрибутов не обязательно должны запоминаться вместе с ключами таким способом, как это показано ниже.

| Дело_№ | Фамилия Имя Отчество    | Пол | Звание | Рожд       | Отдел | Должность           |
|--------|-------------------------|-----|--------|------------|-------|---------------------|
| 58231  | Сидоров Иван Петрович   | муж | 06     | 1971.03.12 | 002   | Начальник станции   |
| 12874  | Смирнова Анна Ивановна  | жен | 08     | 1988.10.11 | 014   | Оператор            |
| 31774  | Шмутькин Ганс Абрамович | муж | 11     | 1999.08.18 | 002   | Инженер-программист |

| Звание |                      |
|--------|----------------------|
| ID     | Name                 |
| 04     | Полковник            |
| 05     | Подполковник         |
| ...    |                      |
| 17     | Кадет                |
| 18     | Вольноопределяющийся |

| Должность |                     |
|-----------|---------------------|
| 08        | Начальник станции   |
| 09        | Инженер-программист |
| 10        | Оператор            |
| ...       |                     |
| 57        | Вахтер              |
| 58        | Повар               |

Почти всегда наблюдается существенная избыточность значений атрибутов в тех случаях, когда они хранятся так, как показано в атрибуте **Должность**. Для того чтобы устранить избыточность, значения атрибутов могут храниться отдельно и снабжаться указателями на них со стороны ключей.

## Вторичные ключи

Можно использовать ключ, который идентифицирует не уникальную запись или кортеж, а все записи или группы, имеющие определенное свойство. Такой ключ называется *вторичным*.

Значение атрибута **ЦВЕТ**

Как вторичный ключ может быть использован атрибут **Отдел**. Этот ключ может быть использован для идентификации тех служащих, которые числятся в определенном отделе.

Иногда таблица имеет много вторичных ключей, которые используются для поиска записей с данными характеристиками.

Связь вторичного ключа с элементами данных или с группами, к которым он относится, может быть реализована различными способами. Один из таких способов – использование вторичного индекса. Вторичный индекс использует вторичный ключ как вход, а на выходе предоставляет первичный ключ, в результате чего может быть идентифицирована нужная запись или группа записей.

Элементарная форма вторичного индекса – инвертированный список.

| Звание | Дело_№ |
|--------|--------|
| 05     | 13875  |
| 05     | 76014  |
| 06     | 58231  |
| 06     | 32960  |
| 08     | 12874  |

| Должность | Дело_№ |
|-----------|--------|
| 08        | 58231  |
| 08        | 32960  |
| 09        | 31774  |
| 09        | 11749  |
| 09        | 55813  |

Инвертированный список содержит все значения вторичного ключа и хранит вместе с каждым его значением соответствующие идентификаторы записи.



## Инвертированные файлы

Существуют два основных способа, с помощью которых данные могут быть организованы и использованы.

Первый способ определяется тем, что каждый кортеж содержит значения атрибутов данного объекта или N:1 ссылки на справочные таблицы, содержащие атрибуты (LUT– LookUp Table).

Второй способ является инверсией первого. С помощью этого способа могут быть получены идентификаторы объектов, связанных с данным атрибутом.

Первый способ хранения данных полезен для ответа на вопрос: *каковы свойства данного объекта?*

Второй – для ответа на вопрос: *какие объекты имеют данное свойство?*

*Полностью инвертированный файл* – это такой файл, который хранит идентификаторы объектов, связанные с конкретным значением каждого атрибута.

*Частично инвертированный файл* является более простым и хранит идентификаторы объектов, связанные со значениями некоторых (но не всех) атрибутов.

## Предикаты

*предикат* – это простое условие (терм) или логическая комбинация простых условий;

*простое условие* – это операция сравнения двух выражений;

*выражение* состоит из операций с константами и именами полей;

*константа* – это значение predeterminedного типа, например целое число или строка.

## Запросы

| Форма   | Тип запроса  | Пример   |
|---|--|--|
| $A(E) = ?$  | Обычный запрос атрибута.<br>Каково значение атрибута $A$ объекта $E$ ?   | Заработок торгового агента № 271 за последний месяц  |
| $A(?) \text{ Rel } V$<br>Rel: $\{=, \neq, <, >, \geq, \leq, \equiv\}$ | Какие объекты имеют заданное значение атрибута.<br>Запрос в инвертированный файл: Какой объект $E$ имеет значение атрибута $A$ , равное (неравное, меньше, больше, ...) $V$ ?        | Кто из торговых агентов заработал больше 2000 долл. за последний месяц?                    |
| $?(E) \text{ Rel } V$   | Перечислить все атрибуты, имеющие заданный набор значений для данного объекта.<br>Запрос менее простой: $?(E)=V$ — Какой атрибут или какие атрибуты объекта $E$ имеют значение $V$ ? | За какие месяцы заработки торгового агента No 271 превысили 2000 долл?                     |
| $?(E)=?$  | Запрос на получение всей информации о данном объекте. Запрос на значения всех атрибутов объекта $E$ .  | Сообщить всю хранимую информацию о торговом агенте No 271                                  |
| $A(?)=?$  | Перечислить значения данного атрибута для каждого объекта. Запрос на значения атрибута $A$ для всех объектов.  | Перечислить заработки за последний месяц каждого торгового агента                          |
| $?(?) \text{ Rel } V$   | Перечислить все атрибуты объектов, имеющие данное значение.<br>Сложный запрос на все атрибуты всех объектов, имеющие значение $V$ .  | Для каждого торгового агента определить месяц, когда его заработок превышал 2000 долларов. |

