

# Модели данных

## А5\_Многомерные модели данных



Московский государственный технический университет  
имени Н.Э. Баумана

**Факультет ИБМ**

Июль 2024 года

Москва

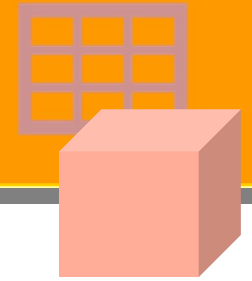
Артемьев Валерий Иванович © 2024

# Модели реляционных витрин данных ROLAP

- Определение витрины данных
- Многомерная модель данных (гиперкуб)
- Основные элементы гиперкуба (факты, измерения, атрибуты)
- Простые и иерархические измерения (сбалансированные, неровные, несбалансированные, альтернативные иерархии)
- Таксономия многомерной модели данных
- Пример таксономии гиперкуба
- Схема «снежинка» (snowflake)
- Схемы «звезда» (star) и «созвездие» (constellation)

Актуализировать,  
когда устоится

# Витрина данных (Data Mart)



*Тематическая база данных, содержащая информацию по отдельным аспектам деятельности организации, предназначенная для обработки средствами бизнес-аналитики.*

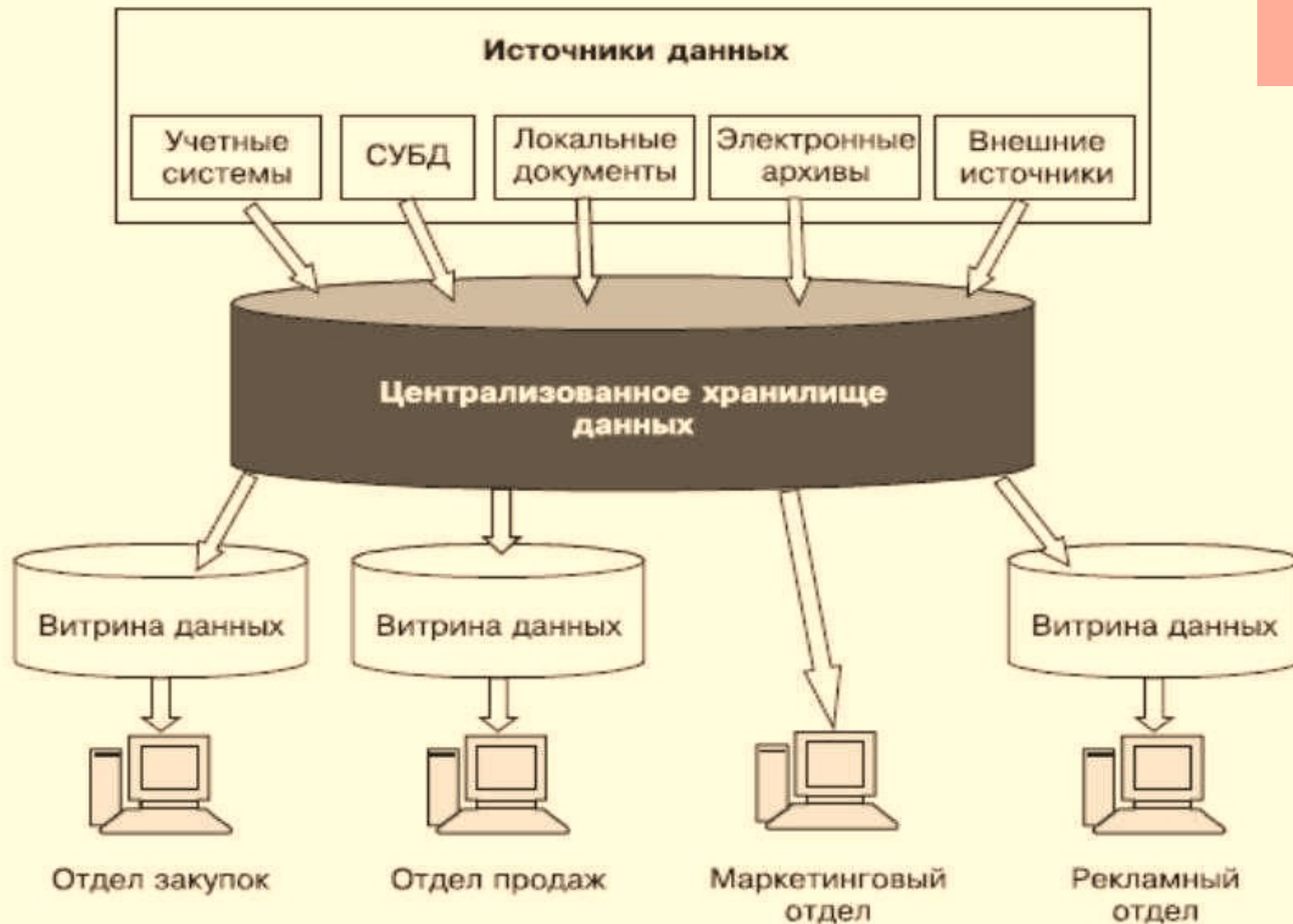
Основные цели выделения витрин:

- приближение данных к конечным пользователям,
- ограничение только необходимыми данными,
- повышение оперативности путём агрегирования данных,
- ориентация структуры данных на бизнес-аналитику.

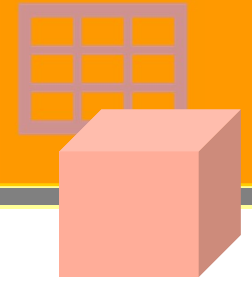
Могут быть независимыми от хранилища данных или зависеть и быть частью хранилища данных.

Часто используется многомерная абстракция для моделирования данных, реализации баз данных и обработки данных.

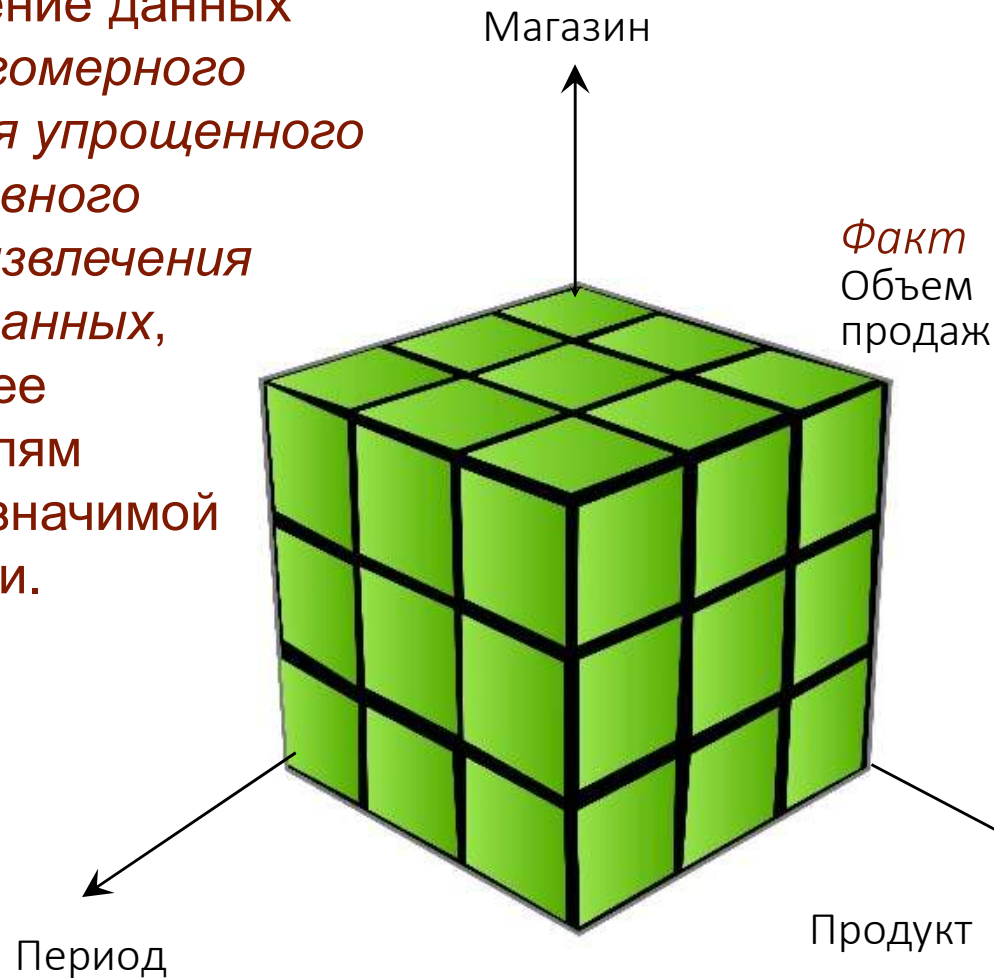
# Место витрин данных в аналитических системах



# Многомерная модель данных (гиперкуб, факты, измерения)

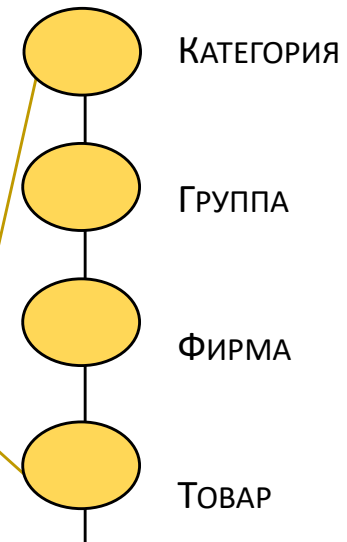


Представление данных в виде *многомерного массива* для упрощенного и эффективного хранения, извлечения и анализа данных, облегчающее пользователям получение значимой информации.



*Гиперкуб* – многомерный массив данных, содержащий факты в нескольких измерениях.

*Уровни*



Измерения могут быть *простыми*, а также *иерархическими*

**Обобщение реляционной модели от таблиц к массивам**

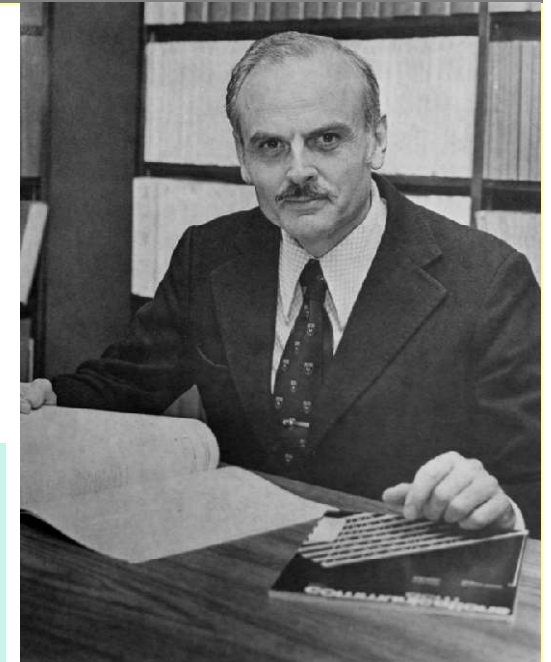


# Интерактивная аналитическая обработка (On Line Analytical Processing, OLAP)

*Технология анализа данных в различных разрезах и с разной степенью детальности, осуществляемого бизнес-пользователями в интерактивном режиме в терминах своей предметной области.*

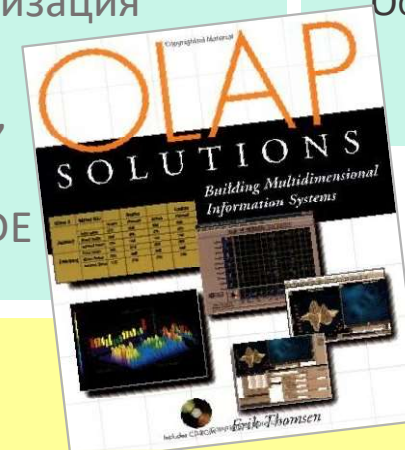


- Описательный и разведочный анализ  
Многомерный анализ данных (детализация и укрупнение)
- Самообслуживание пользователей в бизнес-терминах  
Привлечение опыта и интуиции пользователя
- «Ручная раскопка» данных (детализация и укрупнение)  
Интерактивные запросы и отчёты, расчёты «на лету»  
Инструменты NO CODE / LOW CODE

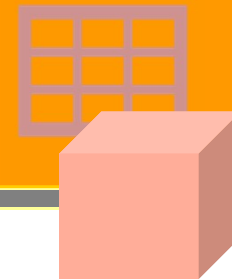


EDWARD CODD  
IBM 1993

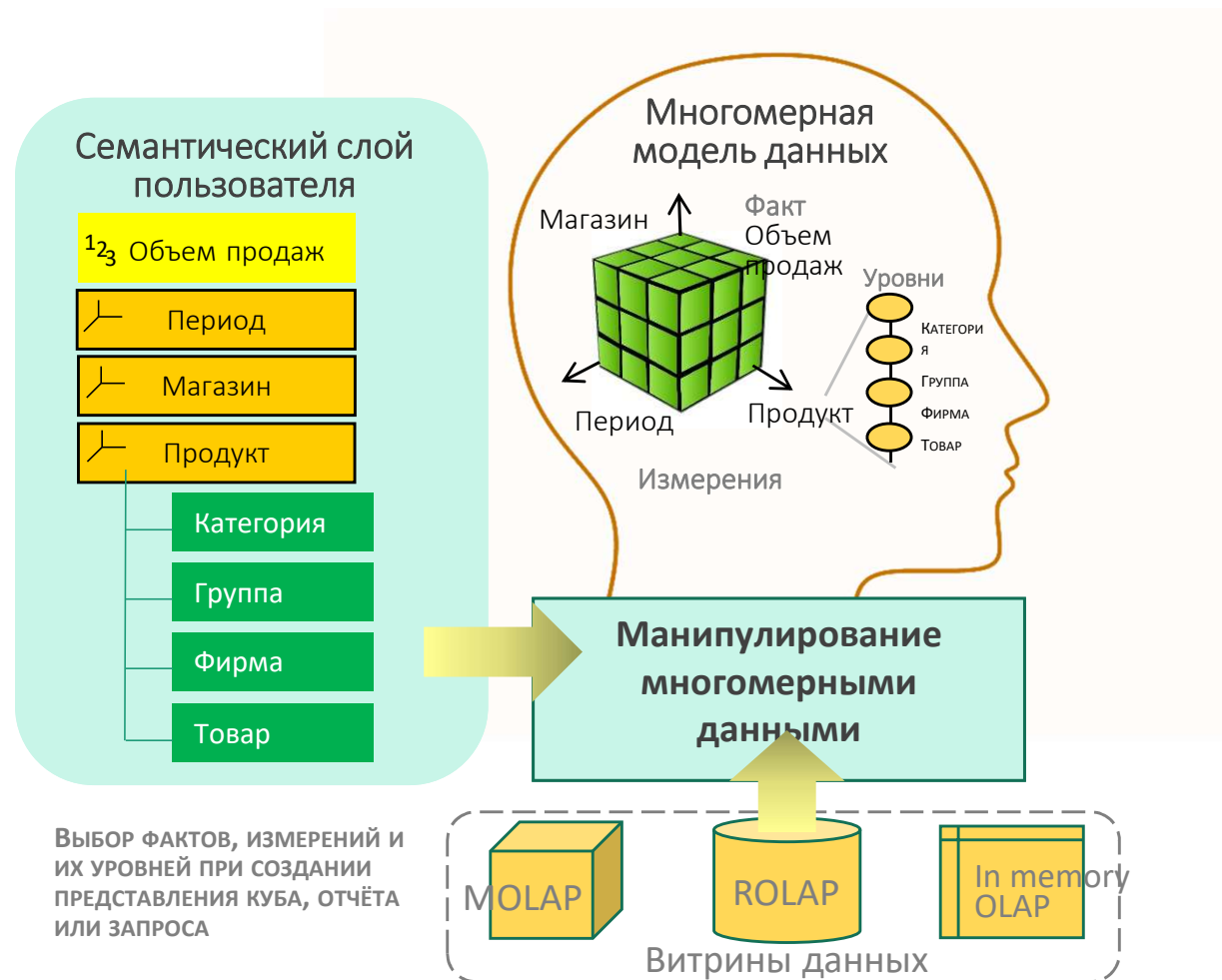
ОСНОВАТЕЛЬ РЕЛЯЦИОННОЙ  
МОДЕЛИ ДАННЫХ



# Многомерные свойства OLAP

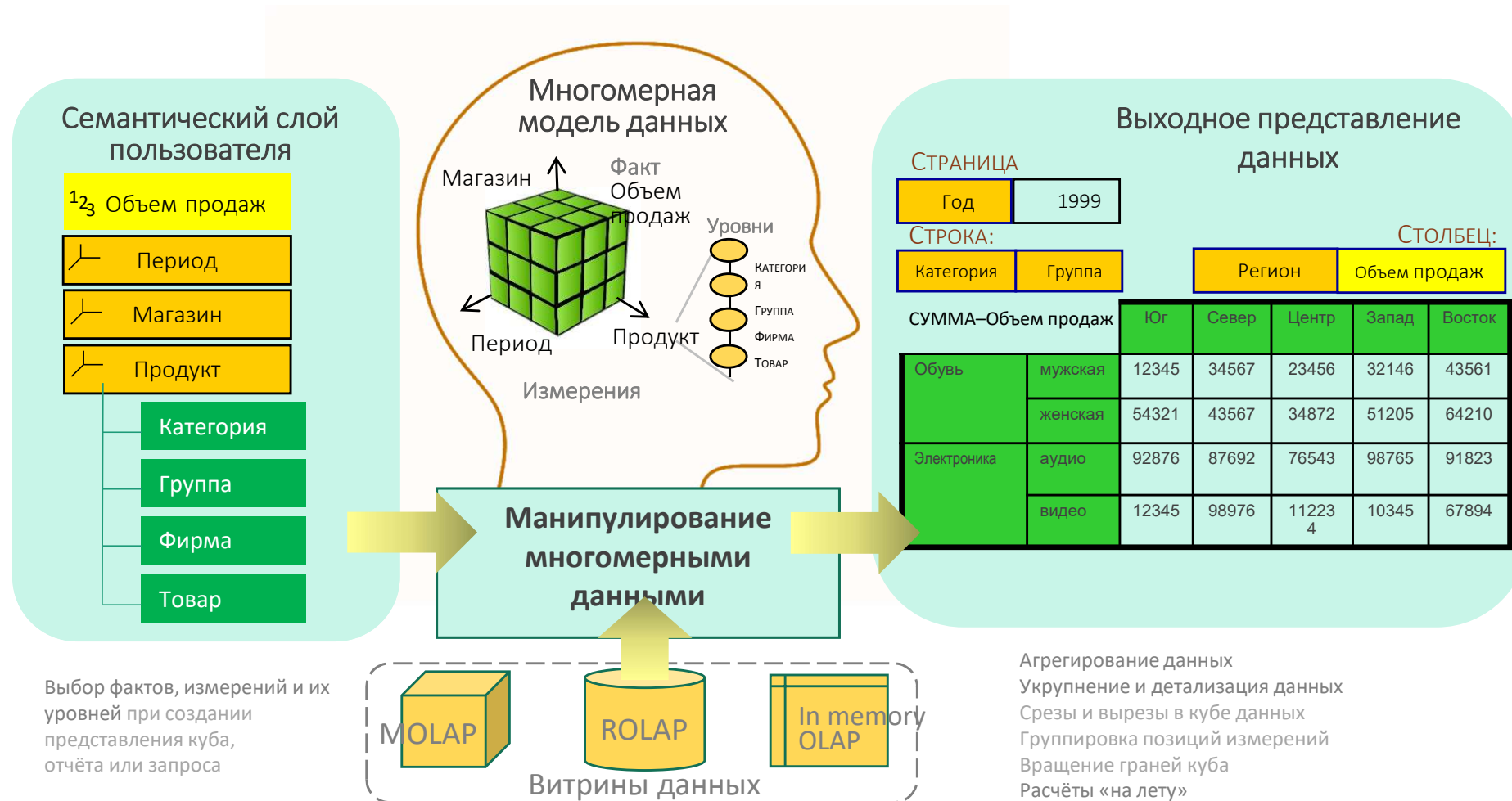
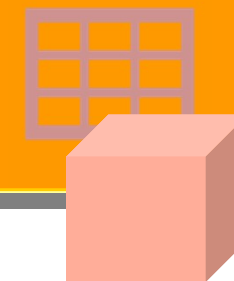


# Многомерные свойства OLAP

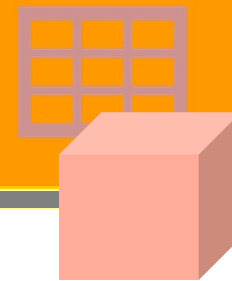




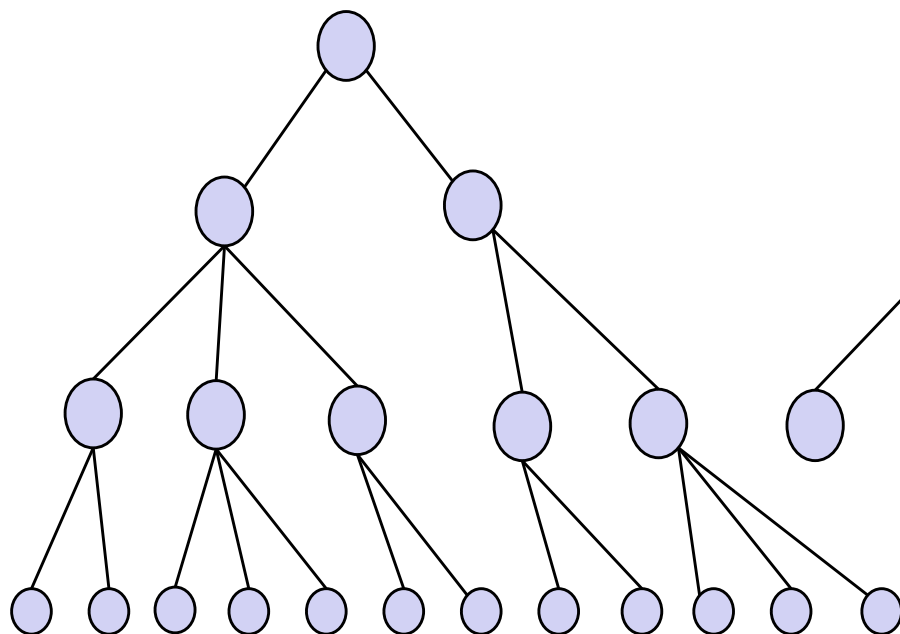
# Многомерные свойства OLAP



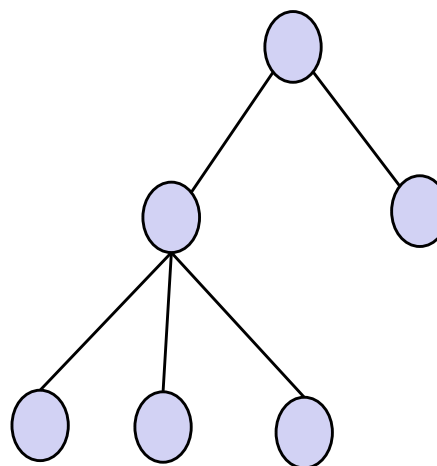
# Виды иерархий измерений



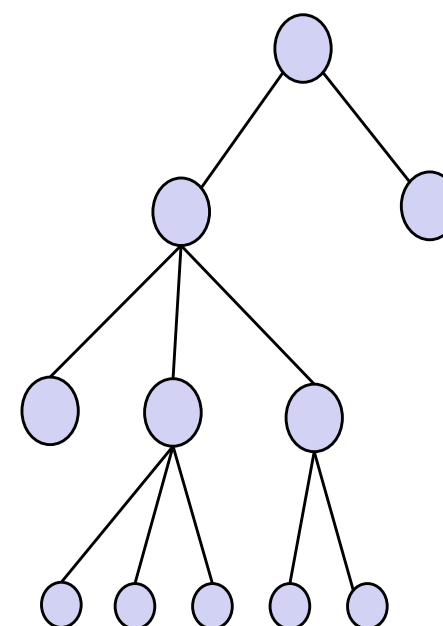
Сбалансированная



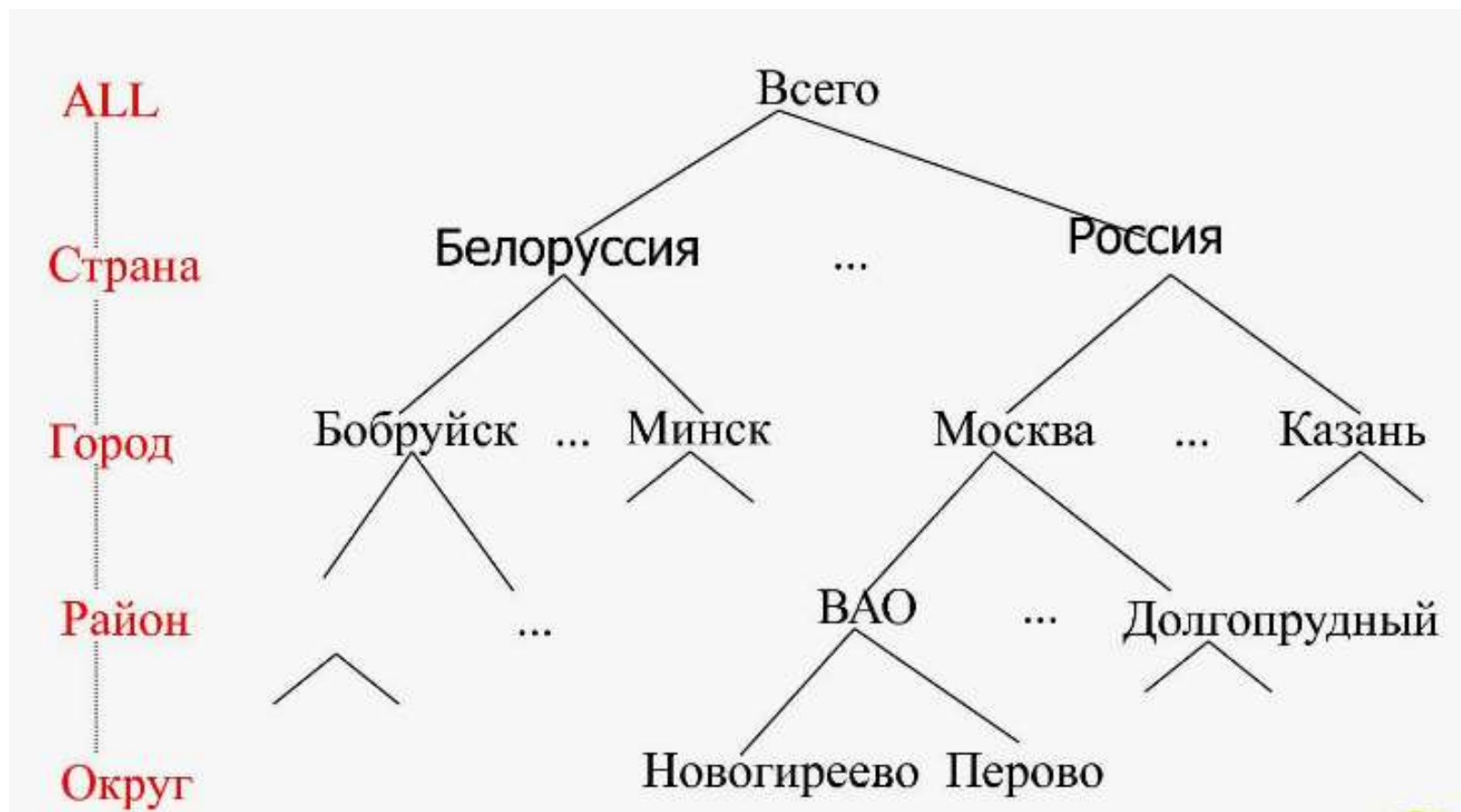
Неровная



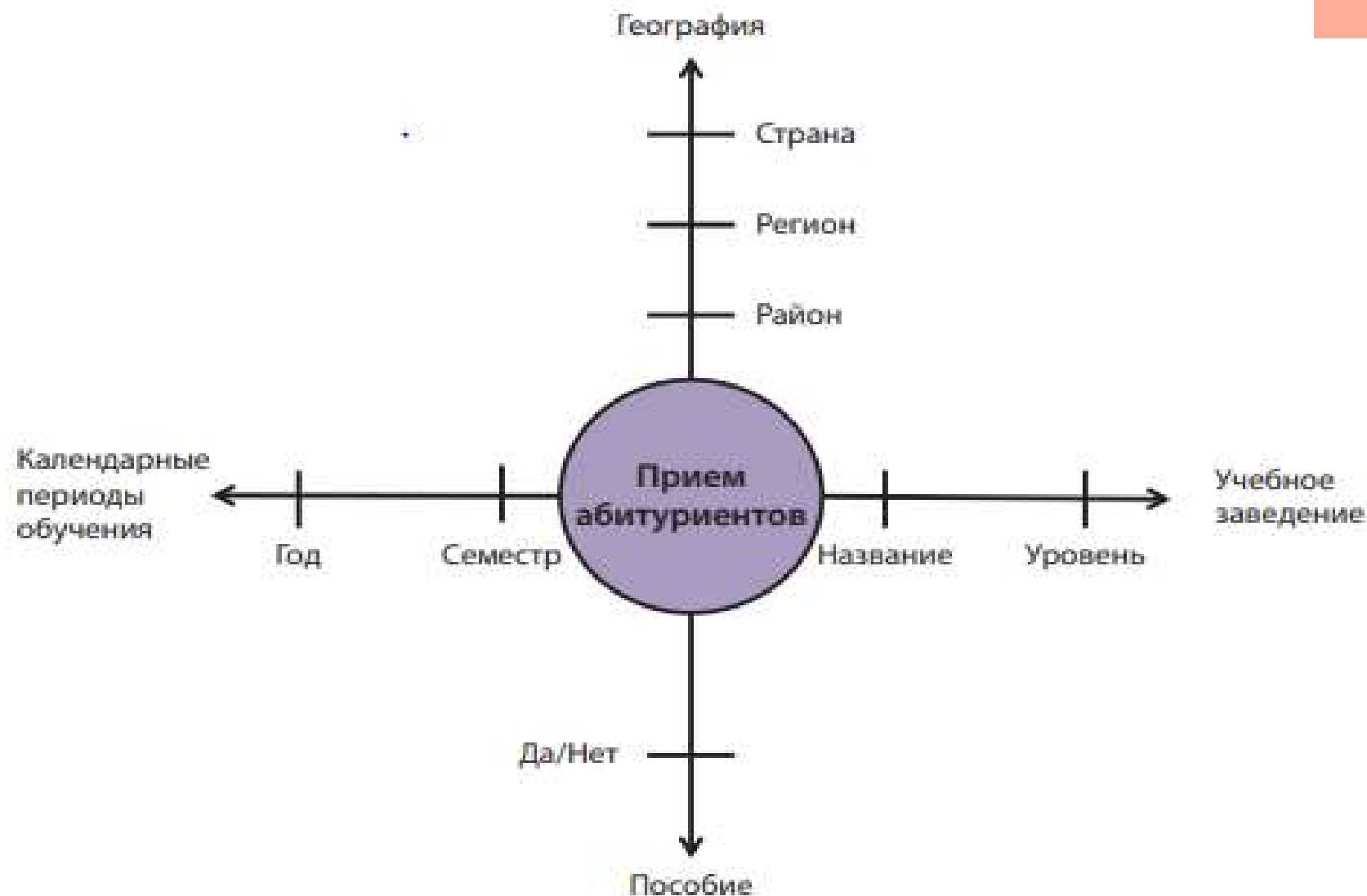
Несбалансированная



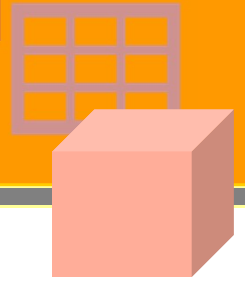
# Пример иерархического измерения



# Осевая нотация для представления многомерной модели данных



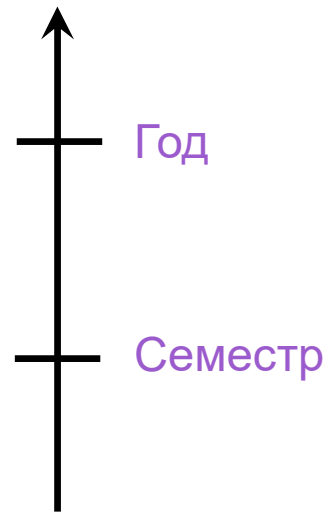
# Другое представление многомерной модели данных



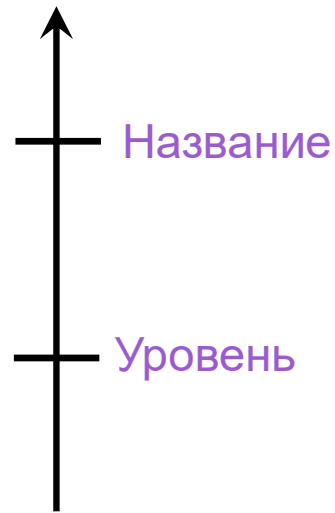
Переменная



Календарь



Учебное  
заведение



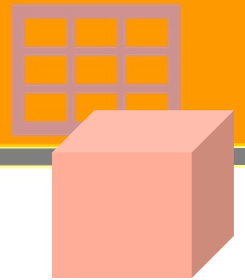
География



Пособие



# Таксономия многомерной модели данных (пример)



## ⊕ Приём абитуриентов (гиперкуб)

### ⊕ ФАКТЫ

⊖ Зачислено

### ⊕ ИЗМЕРЕНИЯ

#### ⊕ Календарь

⊖ Год

⊖ Семестр

#### ⊕ Учебное заведение

⊖ Название

⊖ Уровень

#### ⊕ География

⊖ Страна

⊖ Регион

⊖ Район/



# Реализация многомерных данных



**MOLAP** – многомерная база данных – многомерный массив на диске с прямо адресуемыми ячейками. Агрегаты и детали могут храниться в одном гиперкубе. Как правило, урезанный атрибутный состав.

**ROLAP** – реляционная база данных, имитирующая гиперкуб. Агрегаты хранятся отдельно от деталей. Широкий атрибутный состав.

**In memory OLAP** – многомерная база данных в основной памяти, обычно в сжатом виде с хранением данных по колонкам.



# Структуры данных ROLAP

В реляционном представлении многомерные данные организованы в таблицы двух видов:

**1. Таблицы фактов** содержат количественные данные для дальнейшего анализа:

- *транзакционные факты* на основе отдельных событий;
- «*моментальные снимки*» (snapshot) на основе состояний объекта в определённые моменты времени;
- *факты, связанные с элементами документа*;
- *факты, связанные с событиями и состоянием объекта без подробностей*.

**2. Таблицы измерений** предоставляют контекстную и описательную информацию о данных в таблицах фактов. Часто представляется в виде *схемы снежинки* или *звезды*.

# Структура таблицы фактов ROLAP

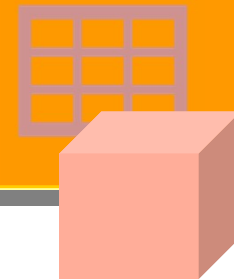


Таблица фактов содержит:

- факты одной степени детальности, соответствующие нижнему уровню иерархии
- внешние ключи таблиц измерений, которые являются составным первичным ключом.

## ПРОДАЖИ

(PK, FK) Календарь\_Ид  
(PK, FK) Магазин\_Ид  
(PK, FK) Продукт\_Ид  
Сумма  
Количество

Агрегаты вычисляются «на лету» по первичным фактам или хранятся в отдельных таблицах фактов.

# Структура таблицы измерений



Таблица измерений содержит:

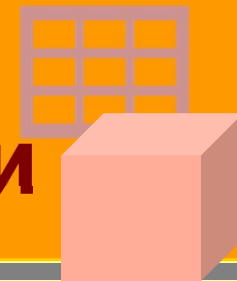
- по одной записи для каждого члена нижнего уровня
- первичный ключ записи
- характеристики нижнего уровня детальности
- квалификаторы уровней агрегирования
- все характеристики верхних уровней иерархии.

## **КАЛЕНДАРЬ**

(PK) Месяц\_Ид  
Месяц\_имя  
Месяц\_номер  
Квартал\_обозн  
Квартал\_номер  
Год

В таблице измерений не рекомендуется хранить смесь разных степеней детальности, это затрудняет обработку и чревато ошибками.

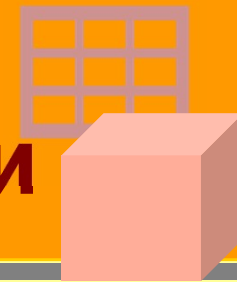
# Медленно меняющиеся размерности



## Slowly Changed Dimension (SCD)



# Медленно меняющиеся размерности

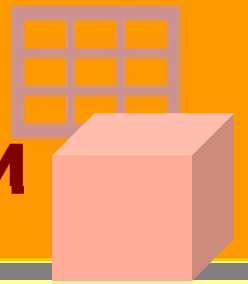


## Slowly Changed Dimension (SCD)





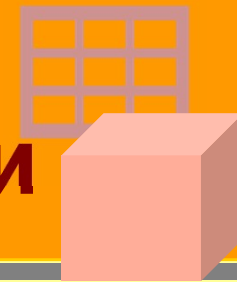
# Медленно меняющиеся размерности



## Slowly Changed Dimension (SCD)



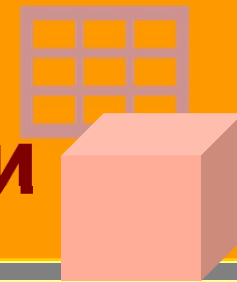
# Медленно меняющиеся размерности



## Slowly Changed Dimension (SCD)



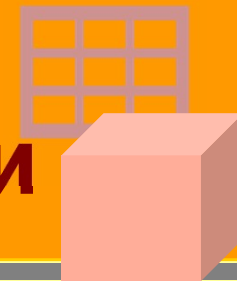
# Медленно меняющиеся размерности



## Slowly Changed Dimension (SCD)



# Медленно меняющиеся размерности



## Slowly Changed Dimension (SCD)



# Медленно меняющиеся размерности

## Slowly Changed Dimension (SCD)



# Медленно меняющиеся размерности



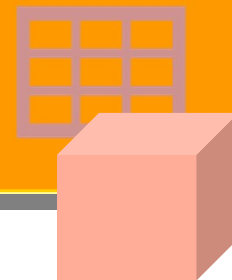
## Slowly Changed Dimension (SCD)



Может использоваться комбинация типов.



# Служебные атрибуты историчности



Для часто используемого способа SCD2 используются дополнительные атрибуты:

- **Effective\_Date** – дата начала действия записи
- **End\_Date** – дата конца действия записи
- **Actual\_Flag** – состояние актуальности

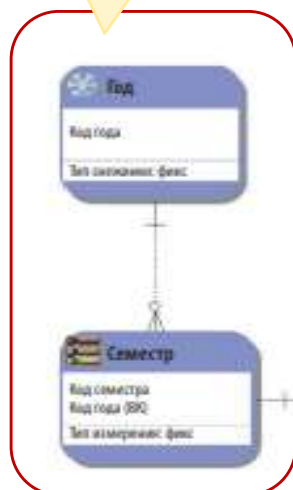
Для происхождения загруженных данных используются дополнительные атрибуты:

- **Load\_Date\_Time** – дата и время загрузки данных
- **Source\_Code** – код источника данных

# Схема «снежинка» (snowflake)



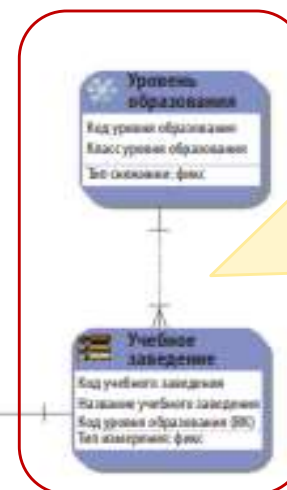
Иерархическая  
размерность  
**Календарь**  
из 2-х уровней  
(справочников)



Иерархическая  
размерность  
**География**  
из 3-х уровней  
(справочников)



Иерархическая  
размерность  
**Образование**  
из 2-х уровней  
(справочников)

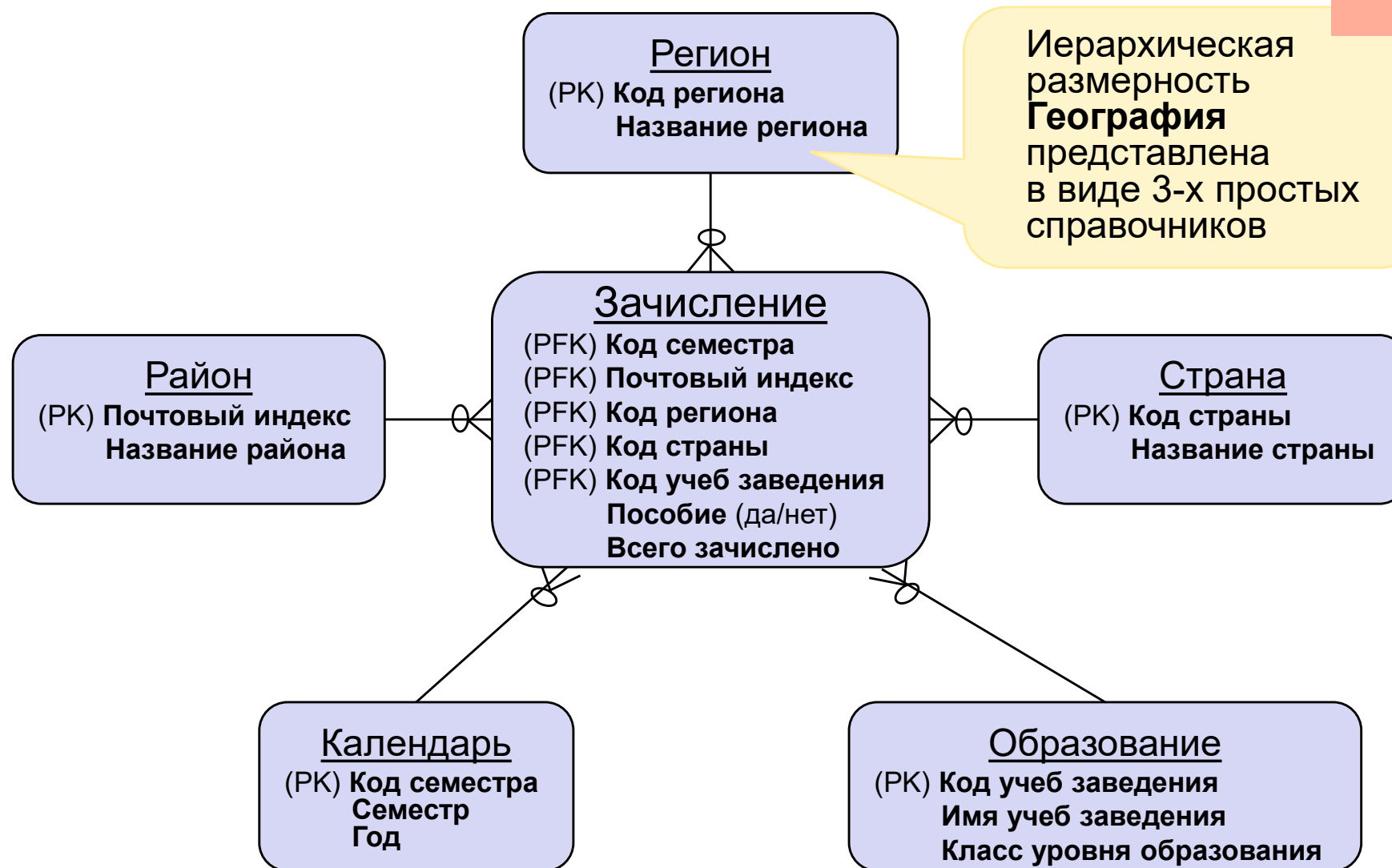
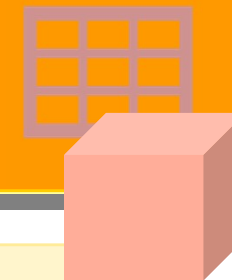


**Таблица фактов**  
имеет 5 атрибутов  
и 3 соединения



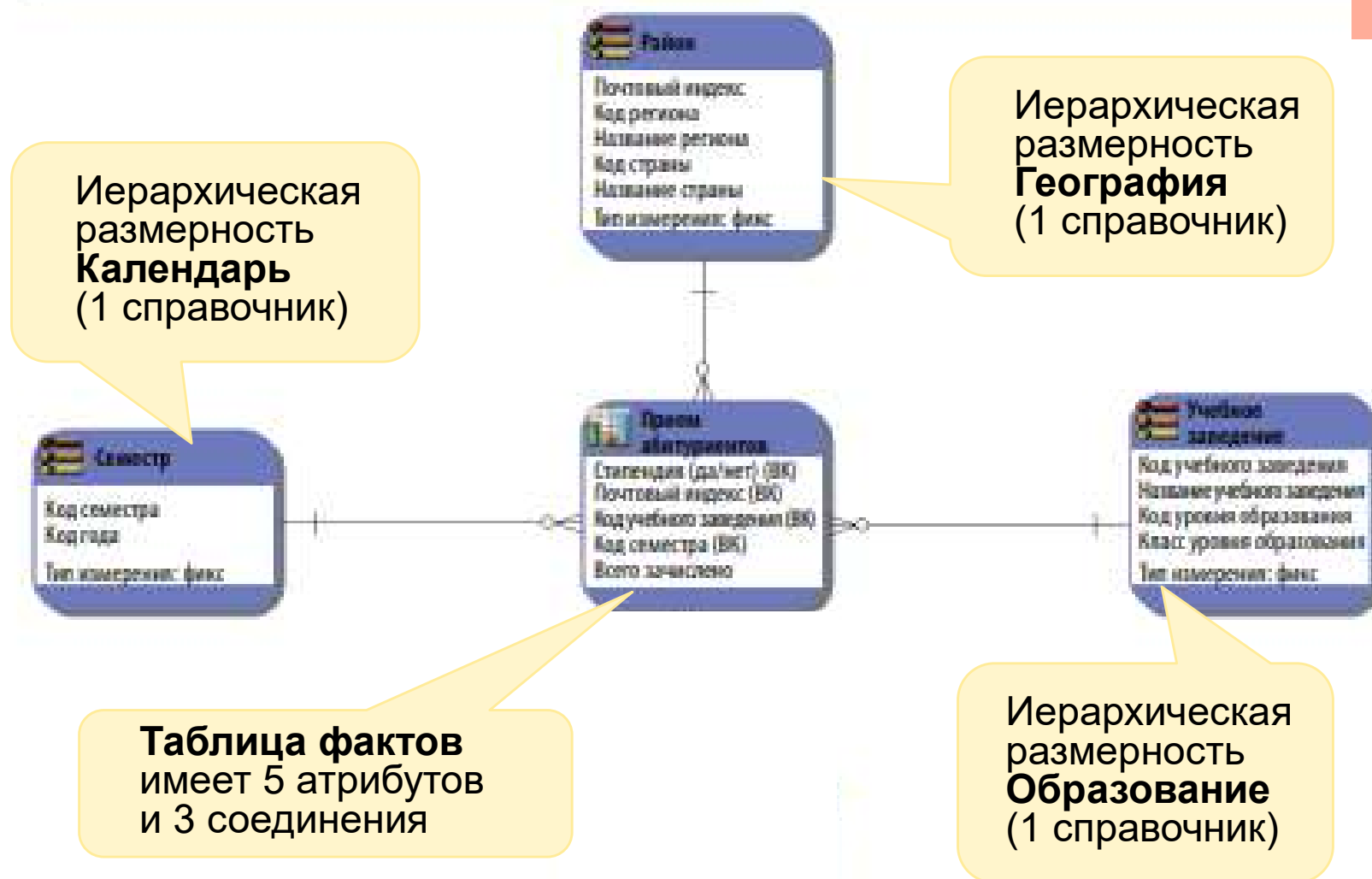
**Высокая нормализация – всего 7 соединений таблиц,  
3 каскадных соединения – снижение производительности**

# Комбинированная схема ROLAP



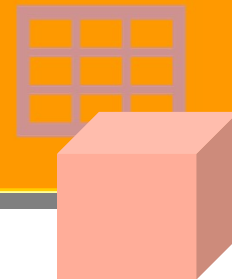
Комбинированное решение – 5 соединений без каскадов, но рост объёма хранения (7 полей вместо 5)

# Схемы «звезда» (Star Schema)



Дополнительная денормализация – 3 соединения без каскадов

# Максимально денормализованная схема ROLAP



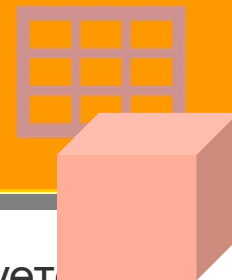
**Единственная плоская таблица.**

Объём хранения определяется количеством атрибутов в таблице фактов

**Таблица фактов**  
имеет 11 атрибутов  
и не имеет  
соединений

**Зачисление**  
(PK) ID  
Семестр  
Год  
Почтовый индекс  
Название района  
Название региона  
Название страны  
Имя учеб заведения  
Класс уровня  
образования  
Пособие (да/нет)  
Всего зачислено

Объём увеличился более чем в 2 раза (11 против 5 атрибутов),  
что снижает производительность, но нет операций соединения

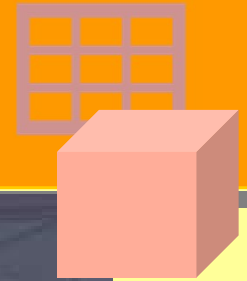


# Вместо заключения

- При анализе данных средствами бизнес-аналитики (BI/OLAP) используется *многомерное представление данных в виде гиперкубов, фактов и размерностей*.
- Бывают *простые и иерархические размерности*.
- Тематические БД для анализа организуют в виде многомерных **MOLAP** или чаще *реляционных витрин данных ROLAP*.
- В ROLAP витринах существуют *таблицы фактов и таблицы измерений*
- Для поддержки истории изменений имеется несколько типов *медленно меняющихся размерностей*, наиболее частый тип **SCD2**.
- Эти витрины создаются путём *денормализации отношений* в виде схем «снежинка», «звезда» или их комбинации, вплоть до плоской таблицы.
- Объём ROLAP-витрины определяется *размером записи таблицы фактов*, грубо числом её атрибутов (полей).
- *Время обработки запросов существенно зависит от количества соединений, в большей степени от наличия каскадов соединений, а также от размера записи таблицы фактов*.
- Решение находится *путём компромисса*.
- Нормализация важна для транзакционных систем, для аналитических систем *важнее денормализация отношений*.



# Спасибо за внимание!



**Терпения и удачи всем, кто связан  
с моделированием данных**

**Валерий Иванович Артемьев**

**МГТУ имени Н.Э. Баумана, каф. ИУ-5**

**Банк России**

**Департамент данных, проектов и процессов**

**Тел.: +7(495) 753-96-25**

**e-mail: [avi@cbr.ru](mailto:avi@cbr.ru)**