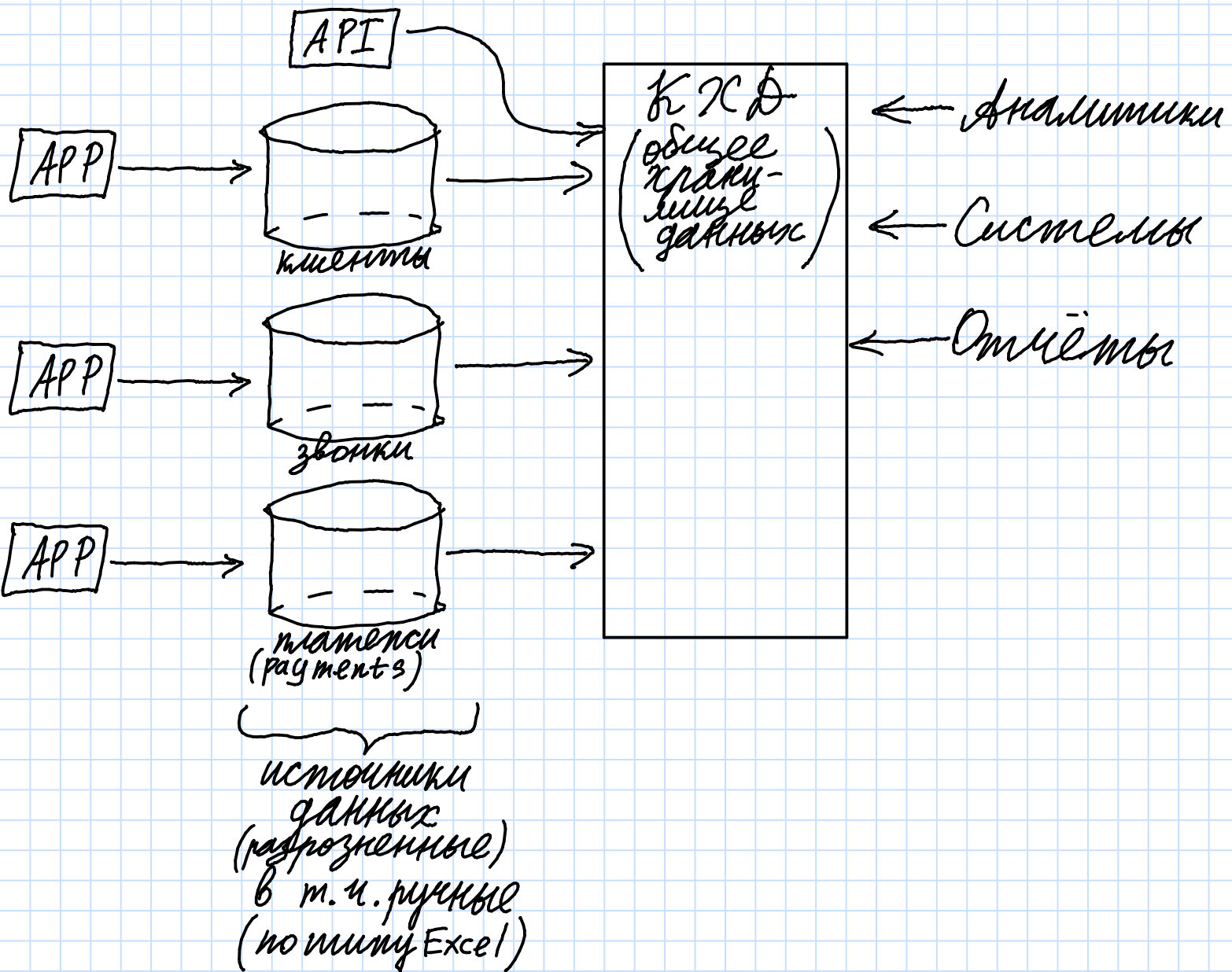Проектирование хранилищ.

K X Д — (корпоративное) хранилище данных
или
D W H — (Data) WareHouses;



Любая передача внутри хранилищ данных — ETL (Extract, Transform, Load);

Способы хранения данных:

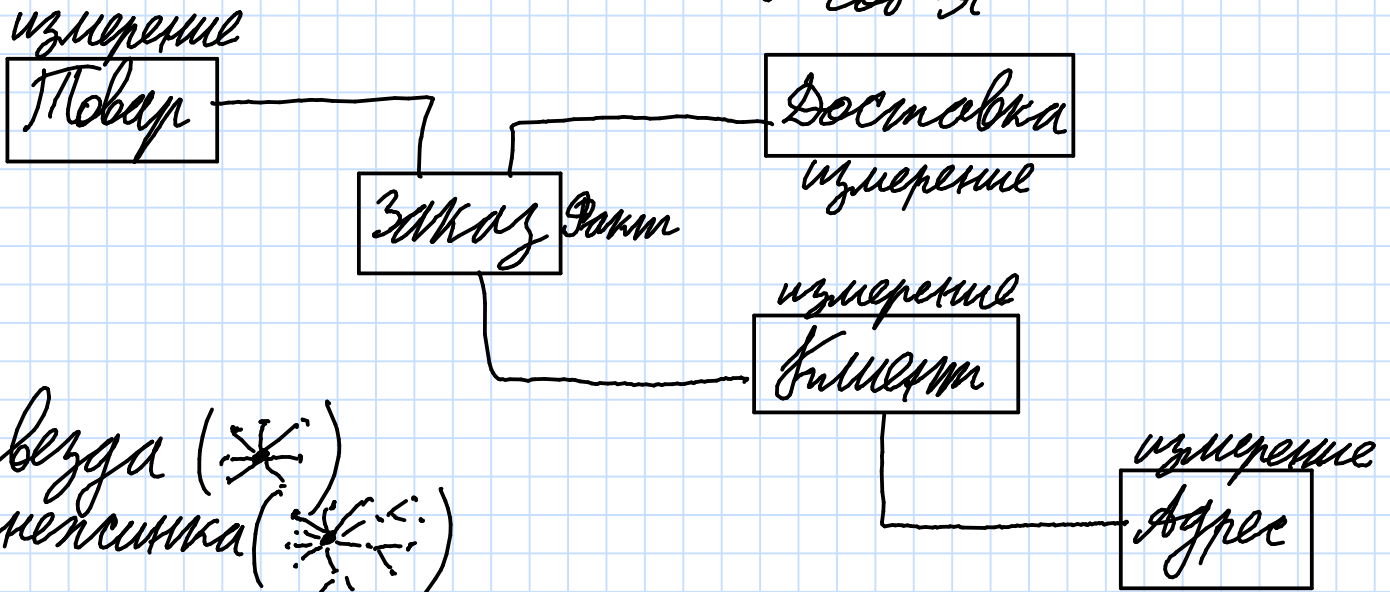Пространственная модель хранения данных.

Оперирует 2 критериями:

— Факт — строка в таблице, кот. явл. точечным

соб-ем (нет продолжи-
тельности)

— Измерение

Все события в системе делятся на эти
2 категории;

__Факт__ — транзакционное соб-е.
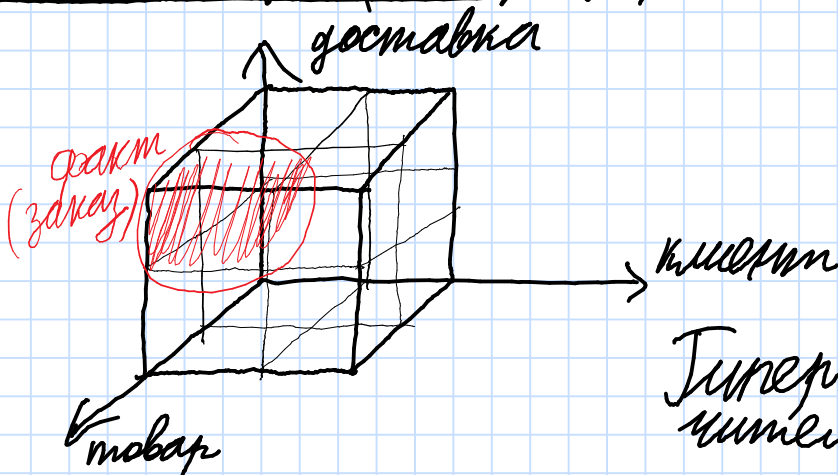↳ у него есть много измерений

__Измерение__ — мера соб-я ( = справочник)

могут иметь
под собой точечные
соб-я



измерение
| Товар |

| Заказ Факт |

| Доставка |
измерение

измерение
| Клиент |

измерение
| Адрес |

звезда ( ✳ )
снежинка ( ✳ )

__Измерение__ — сборище звёздочек или снежинок;

__OLAP — куб (гиперкуб)__

доставка

факт
(заказ)

клиент

товар

Excel умеет
визуализировать
гиперкубы!
(отчёты в разных
разрезах)

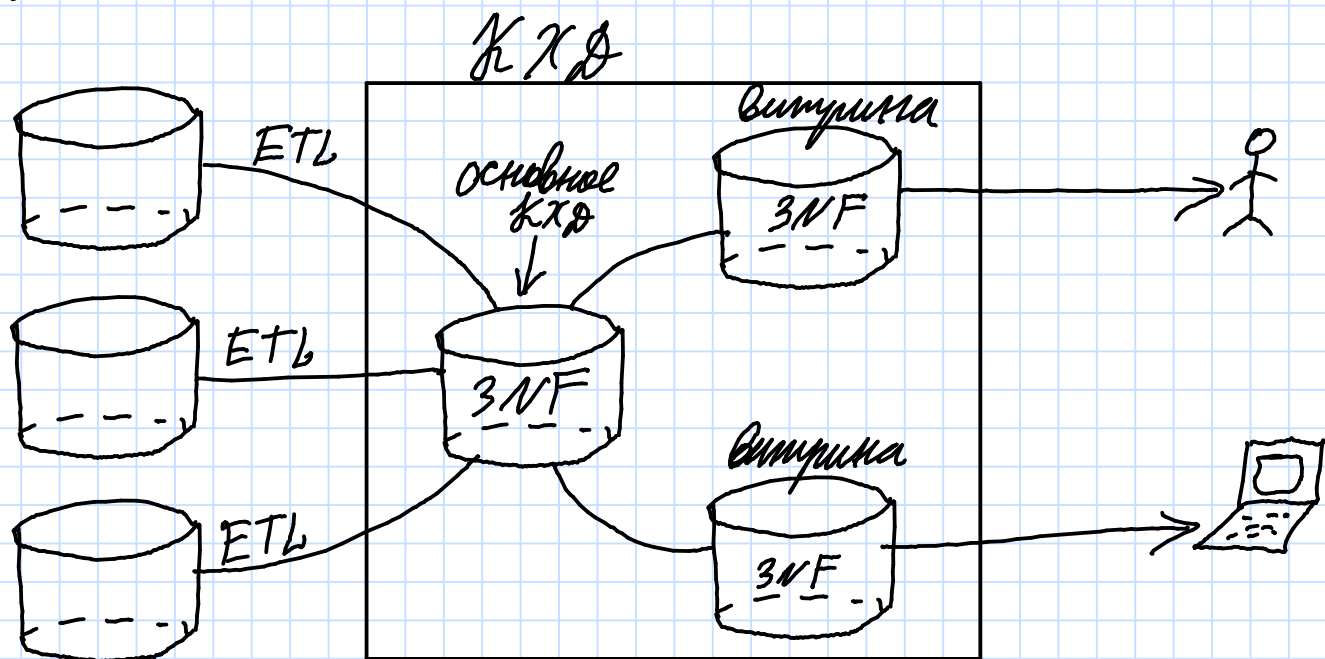Гиперкуб построен исклю-
чительно на звёздочках;

Исторически сложились 2 модели хранилищ данных:

# ① Модель хранилища данных по Инмону:

Данные должны быть структурированы;

КХД — единая (·) „правды";

Данные должны быть

— чистыми;
— достоверными;
— непротиворечивыми;
— подчиняться 3NF;

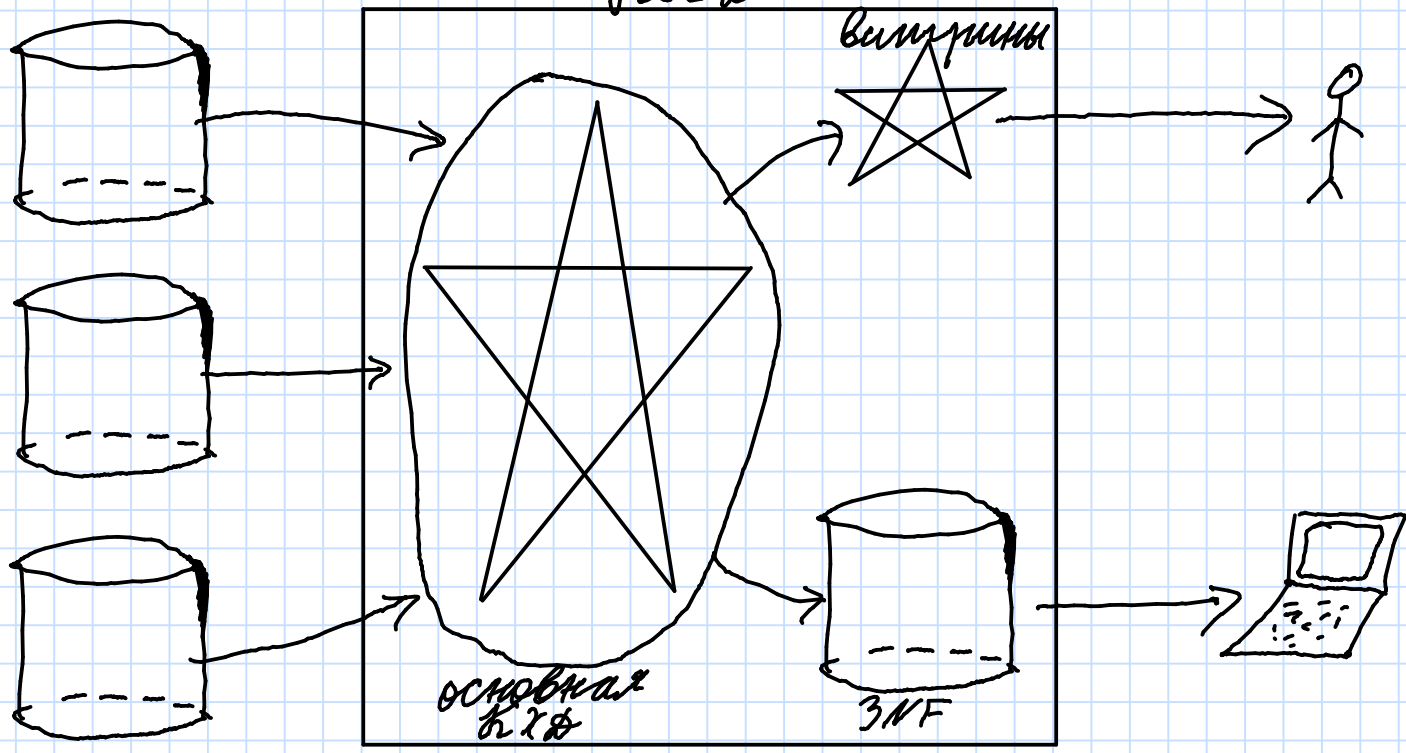3NF должна лежать в основе всех хранилищ

3NF:

$\oplus$ Чистота данных( все данные красивые);

$\oplus$ V хранилища (данные неизбыточные);

$\oplus$ Гибкость ( относительно легко добавлять данные);

$\ominus$ много join'ов (очень много времени);

Колоночные MPP системы очень тяжело джойнятся (по типу Klickhouse)

$\ominus$ время начальной разработки (2-3 мес. вместо 3 часов);

② Модель хранилища данных по Кимбаллу:

КХД

Витрины

основная КХД

3NF

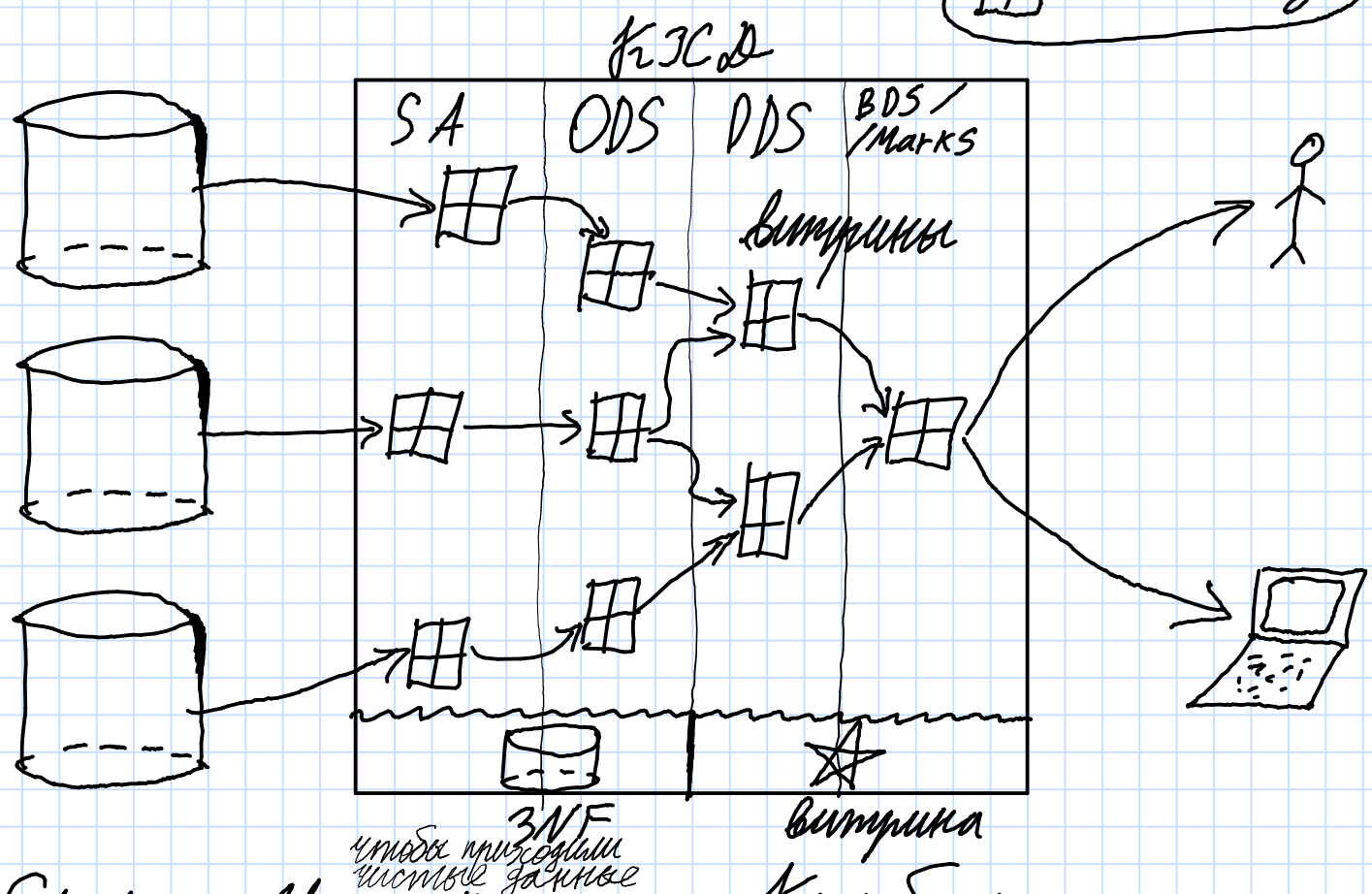Здесь чаще всего не 3NF, а декомпозиция

→Она предполагает предварительную нормализацию ⟹ время разработки не сильно уменьшится;

⊕ время запросов

⊕ звёздочка выглядит приятнее (с т.з. пользователя), чем 3NF;

⊖ Аномалии (данные не очень чистые);

⊖ V хранилища (избыточность данных);

Обе методологии не используются ⟹
⟹ гибридный подход;

⊞ – таблица

КХД



SA    ODS    DDS    BDS/Marks

витрины

3NF
чтобы приходили
чистые данные

витрина

Сначала Инмон, потом Кимбалл;
Слои:
1) Staging - предподготовленные данные;
2) ODS - операционные данные (operation data stage);
SA - сырые данные (JSON, инкременты и т. д.);
ODS - по сути те же источники, только у нас
в хранилище; ⟶нужен для выч-я агрегаций;
DDS - Detail Data Stage; BDS-Buisness Data Stage;

# Data Vault

3 сущности:

— <u>Хаб</u> — имеет опис-е к фактору, в кот. чаще всего помещается набор ключей;

Это ключ: натур. + источник + время появления
(суррогат.) ‾‾‾‾‾‾‾‾‾‾‾‾‾‾‾‾‾‾‾
технические поля

( Напоминает усложнённый факт )

— <u>Сателлит</u> — ключ (ссылка на хаб) + атрибуты + ...
+ тех. поля

Сателлит и хабы — one to many;

( Напоминает измерение )

— <u>Ссылка</u> — связь между хабами + тех. поля
Это всегда связь many to many;

При построении бизнес-слоя (BDS) будем использовать эти 3 сущности;

## Якорная модель хранения

4 объекта: (факт свершения соб-я)

— <u>Якорь</u> — суррогат. + технические поля ;
ключ

— <u>атрибуты</u> — сур. ключ. якоря + поля + тех. поля;

— <u>связь</u> — только many to many;

— <u>узел</u> : атрибуты, код (справочник);


В Data Vault упор на натур. ключи,
а Якорная модель — на суррогатные;

<u>Натуральные ключи</u> VS <u>суррогатные ключи</u> :
- ... - очень долгий и неприятный разгон ...
- сур. ключ должен быть создан при помощи
хеш-функции ;