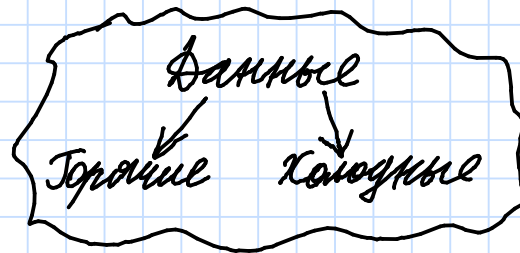Data Lake могут поддерживать только большие IT компании.

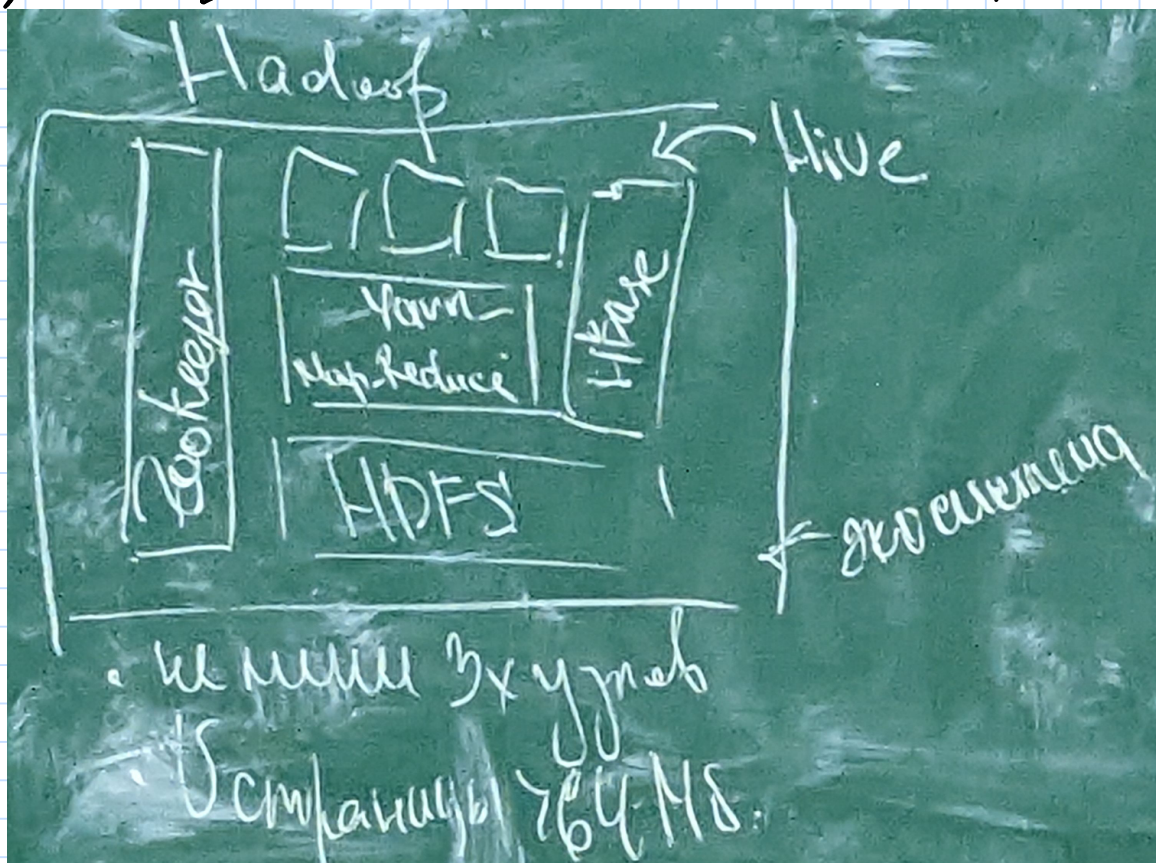---

Hadoop — распределённое хранилище данных (именно холодных) → не часто нужны, но возможно на иметь у себя
(написан на Java)



Данные
↓          ↘
Горячие    Холодные

Тенденция вынесения холодных данных в отдельные сервера ⇒ Hadoop
(ориентировока на большие данные)

Основан на своей файл. системе HDPS
(1 страница — 64 Мб, а не 4/8 Кб)

Компоненты универсальны
(могут использоваться отдельно)
в других системах
        (HDFS исп. в Clickhouse)
Hadoop: хранение и расчёт холодных
данных, фоновые задачи (а ля "запустить
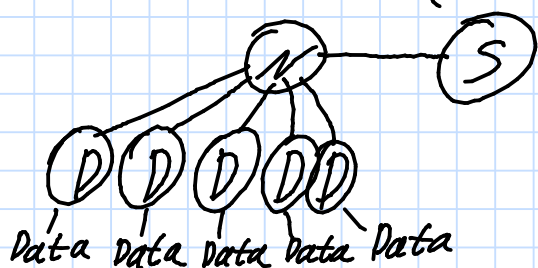                                раз в сутки/
                                /месяц/год")

He realtime!

В крупных компаниях рядом с Hadoop
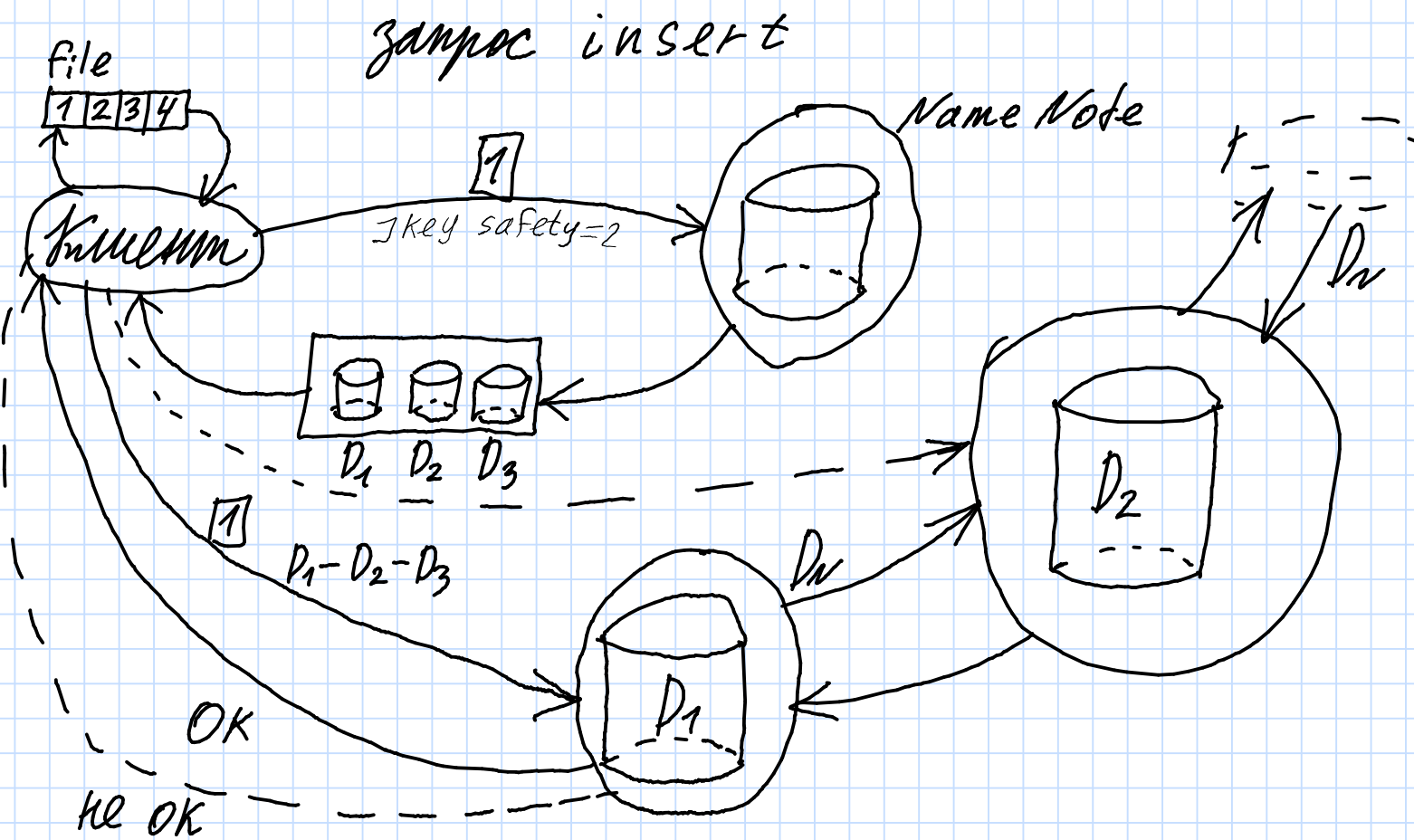стоит OLAP-система (для горячих
данных): Clickhouse, Vertica,...

Аналитические системы на Hadoop
построить же удалось

3 вида узлов:
1) Name Node — для маршрутизации;
( хранит метаданные + контакт с клиентом)

2) Secondary Name Node - дублёр Name Node
+ Выполняет бекапы журналов (репликация)

3) Data Node (желательно ≥ 3 шт.)


Data  Data  Data  Data  Data

Типовое взаимодействие:

запрос insert



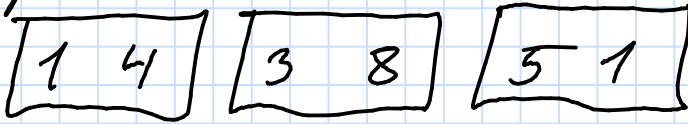Узлы сами сигнализируют о своём состоянии;

Алгоритм Map - Reduce.
Шаги:
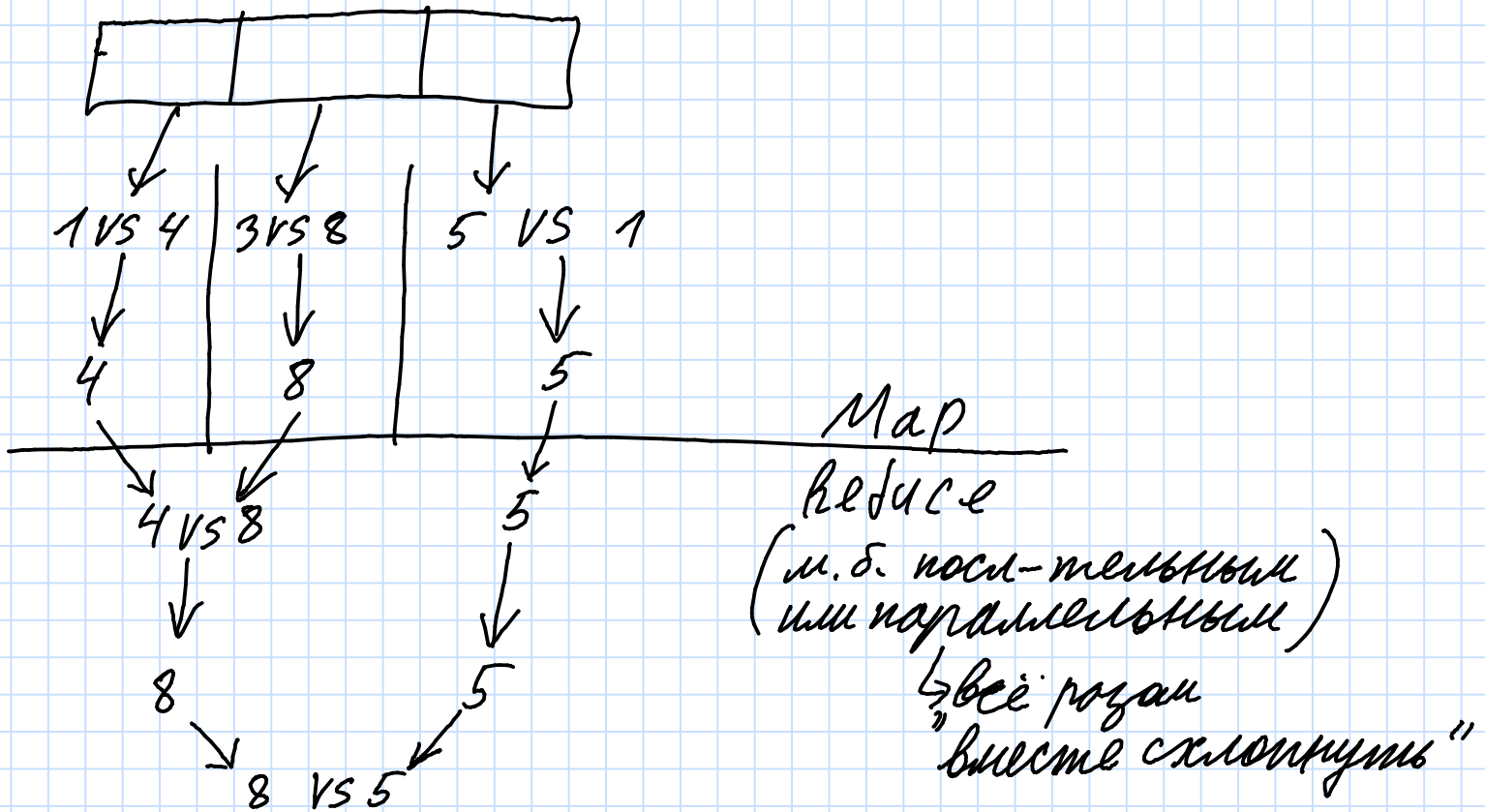1) Map - распределение задач по серверам;
(на потоки)

Hadoop не применяется для посл - тельных и рекуррентных задач, т.к. они плохо распараллеливаются;
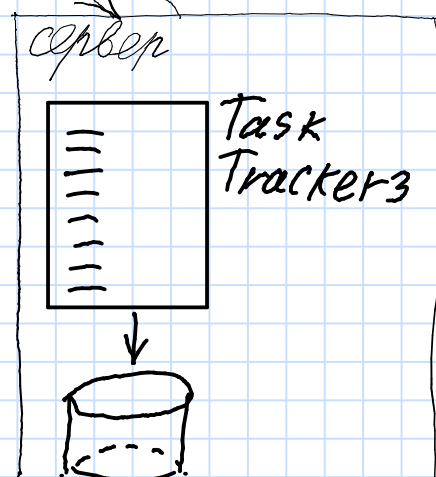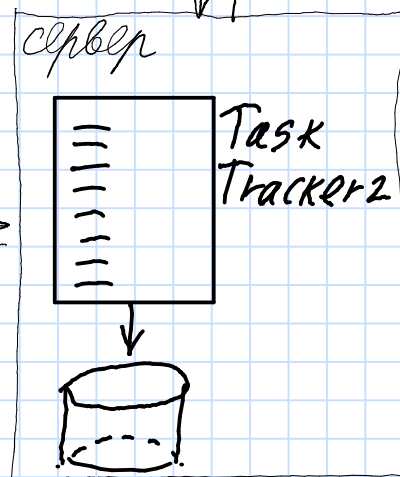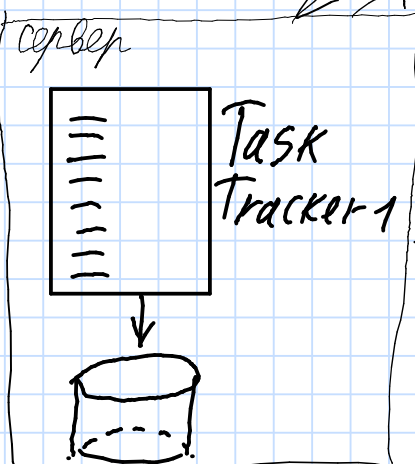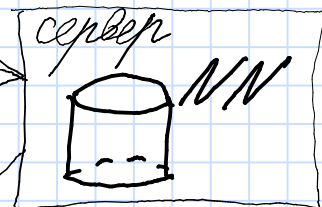
2) Reduce;

# Пример: нахождение максимума в массиве

| 1 | 4 | | 3 | 8 | | 5 | 1 |
|---|---|---|---|---|---|---|---|

max = ?

| | | |
|---|---|---|

1 vs 4      3 vs 8      5 vs 1

4            8            5

───────────────────────────────── Map

Reduce

4 vs 8                   5

8                        5

8 vs 5

(м. б. посл-тельным или параллельным)

↳ «всё разом вместе схлопнуть»

{ NN отвечает за балансировку }



сервер

NN

Клиент

сервер

Job Tracker

map
reduce

сервер

Task Tracker-1

сервер

Task Tracker-2

сервер

Task Tracker-3