

# Data Engineer

Лучшие практики по приготовлению данных. Загрузка, обработка, организация хранения и доступа к данным с использованием современных инструментов

Длительность курса: 122 академических часа

## 1 Data Engineer: задачи, инструменты, умения

### 1 Инженер Данных. Задачи, навыки, инструменты, потребности на рынке.

- Кто этот специалист и зачем?
- Какая ценность для бизнеса?
- Потребность на рынке. Навыки
- Дальнейшее развитие
- Data Driven Decisions
  
- Какие задачи решает?
- Инструменты для решения основных задач
- Введение в большие данные: где используют, экосистема и история развития фреймворков

Домашние задания

- 1 Домашнее задание: анализ рынка Инженер Данных: РФ, USA, EU

Топ потребностей бизнеса, ключевые технологии и умения  
Понять для себя, где и чем бы вам хотелось заниматься  
Фиксация целей на процесс обучения

---

2 **Эволюция подходов работы с данными. Базовые принципы и понятия.**

- CAP theorem, Distributed Computing, MPP (Massive Parallel Processing).
- Эволюция аналитических хранилищ данных
- SQL + Not Only SQL.
- Batch + Stream
- Lambda, Кappa

Домашние задания

1 Строим свою архитектуру

Цель: В данном домашнем задании от участников курса ожидается небольшой архитектурный документ (желательно - не более 3-х страниц). Выберите интересный сервис (Twitter / Uber / LinkedIn / ваша организация или собственный проект) и разработайте для него архитектуру аналитического хранилища данных. Опишите возможные требования к хранилищу, источники и архитектуру хранения. Приведите пример бизнес-кейса, который можно решить с помощью выбранной вами архитектуры. В решении обязательно должна присутствовать архитектурная схема вашего решения, которая должна объяснять откуда к вам поступают данные, как вы планируете их хранить и как вы планируете их отдавать для решения бизнес-кейса. Описание архитектуры - это стандартное задание на архитектурных собеседованиях для Data Engineer. Выполнив данное задание, вы сможете в будущем воспользоваться этими знаниями для того, чтобы качественно и продуманно создавать дизайн для аналитического хранилища данных или понимать как ваше data-driven приложение ложится в экосистему вашей организации. Для рисования схемы советуем использовать бесплатный сервис draw.io

Разместите ответ в виде документа Google Docs и поделитесь ссылкой.

3 **Облачные платформы. GCP, AWS, Azure.**

Облака: Amazon Kinesis, Google Cloud Pub-Sub, Google Dataflow. Cloud functions

Домашние задания

1 Создание пайплайна обработки данных на Google Cloud Platform

Цель: В этом домашнем задании мы создадим простейший пайплайн обработки данных 1) Развернем виртуальную машину, которая будет играть роль продуктового сервиса, генерирующего данные 2) Зальем данные в хранилище GCS 3) Загрузим эти данные в BigQuery для возможности анализа на SQL

Ссылка на инструкцию к домашнему заданию

<https://docs.google.com/document/d/1OENWwlluHYC0tMmxeySGBIFYtISiflbTXXRKSisSuis/edit?usp=sharing>

4 **Дистрибутивы Cloudera и Hortonworks**

- Кто такие Cloudera и HortonWorks и что за экосистемы они строят
- Как собрать кластер Hadoop на основе популярных дистрибутивов CDH и HDP

## 2 Загрузка и форматы данных (Data Ingestion)

- |   |  |  |
|---|--|--|
| 1 | <b>Распределенные файловые системы</b>                           | <ul style="list-style-type: none"><li>- Принципы работы распределенных файловых систем</li><li>- Структура кластера HDFS</li><li>- Тонкости настройки HDFS - конфигурация, защита, обеспечение отказоустойчивости</li></ul> <hr/>  |
| 2 | <b>Инструменты выгрузки данных из сторонних систем - 1 часть</b> | <ul style="list-style-type: none"><li>- Типы систем-источников. Структурированные, полу- и неструктурированные данные. Логи, выгрузки из AC, Clickstream</li><li>- Инструменты для извлечения и загрузки данных - Flume, Sqoop, StreamSets, Fluentd, Debezium, logstash</li><li>- Практические примеры загрузки данных из сервисных баз данных</li></ul> <hr/> |

### 3 Инструменты выгрузки данных из сторонних систем - 2 часть

- Типы систем-источников. Структурированные, полу- и неструктурированные данные. Логи, выгрузки из AC, Clickstream
- Инструменты для извлечения и загрузки данных - Flume, Sqoop, StreamSets, Fluentd
- Практические примеры загрузки данных из сервисных баз данных

#### Домашние задания

##### 1 Создать снейшот аналитической таблицы из операционного хранилища

Цель: Зачастую в Data Lake не требуется хранить всех сырых данных, например в ситуации когда наиболее интересными для пользователя являются какие-либо аналитические выгрузки. В нашем домашнем задании мы напишем пример такого приложения, которое позволяет писать в Data Lake текущий аналитический срез.

1. Скопируйте на свою рабочую станцию (ПК или Google Cloud) данный репозиторий: <https://github.com/renardeinside/vogel>
2. Перейдите в папку проекта (vogel) и запустите все необходимые инстансы командой:  
`docker-compose up --build`
3. Сгенерируйте токен для доступа к Jupyter Notebook (выполните эту команду в проектной папке):  
`docker-compose -f docker-compose.yml exec jupyter-local jupyter notebook list`
4. Перейдите в браузер и введите в нем полученную ссылку с токеном.
5. В Jupyter Notebook перейдите в папку work, в ней откройте ноутбук `postgres_snapshot`.
6. Запустите в нем все шаги до шага Read and join the data.
7. В переменную `query` запишите select-запрос, который выводит следующие данные (доступные таблицы - `customers`, `products`, `orders`):

`id` - id пользователя из таблицы `customers`  
`first_name` - `first_name` пользователя из таблицы `customers`  
`last_name` - `last_name` пользователя из таблицы `customers`  
`total_weight` - суммарный размер всех продуктов данного пользователя (поле `weight` в таблице `products`)  
`load_dttm` - дата/время выполнения команды загрузки

Выполните данную ячейку командой Shift + Enter.

Дополнительные ограничения запроса - необходимо выграть только тех пользователей, у которых `id <= 1005`.  
Информацию по структурам таблиц вы можете прочесть в README.md файле.

8. Проверьте что были записаны верные данные в следующей ячейке.
9. В комментарии к ДЗ пришлите ваш форк данного репозитория с верной sql-командой.

### 4 Форматы данных и их особенности

- Назначение row-based и column-based форматов
- Обзор наиболее распространенных форматов: Avro, Parquet, ORC

## 1 Очереди сообщений

- Kafka, RabbitMQ
- Поточковая обработка (виды обработки, описание Producer–consumer problem, пример архитектурного решения через Kafka, RabbitMQ, NATS)
- Google Dataflow paper (Event time vs processing time и так далее).
- Паттерны stream processing Joins, enricher, router. Event-sourcing.

## Домашние задания

- 1 Архитектурный анализ применимости очереди сообщений для конкретного кейса

Цель: В этом ДЗ мы оцениваем архитектурные решения (очередь сообщений/файловый обмен/RPC) для ваших кейсов, пробуем рассуждать о архитектурных критериях, ключевых требованиях

1. Предложите задачу построения взаимодействия двух приложений близких к вам (примеры - слайды 18, 20)
2. Используя критерии на слайдах 15-16 (или предложите свои) проведите анализ вариантов взаимодействия и сделайте обоснованную рекомендацию

- 2 Анализ стримингового приложения по модели Dataflow

Цель: В этом ДЗ мы изучаем модель Dataflow и рассуждаем о простом стриминг приложении в концепциях Dataflow

1. (слайд 57) Рассмотрим стриминг приложение которое считает остаток на счете дебетовой карты:

- Получает события вида “timestamp, account, amount”, где amount может быть положительным или отрицательным
- Максимальное время задержки доставки событий: 2 секунды
- Если на счете недостаточно средств – приложение выдает алерт и прежний остаток, если достаточно – обновленный остаток

2. Прочитайте Dataflow paper (приложено в материалах) и развернуто ответьте на четыре вопроса описывающее это стриминговое приложение:

- Что считаем? (What?) Что считает приложение? Описание трансформации над событием
- Где по времени события? (Where?) Где во времени (над какими событиями на линии времени) происходят трансформации? Описание используемых окон
- Когда по времени обработки? (When?) Когда по времени обработки материализовать результат трансформации? Описание водяных знаков и триггеров
- Как связаны? (How?) Как связаны последующие результаты трансформации с предыдущими? Описание агрегации (дополнительное задание: опишите как бы вы поступили в случае поздней доставки данных, если задержка более двух секунд)

## 2 DWH. Хранилища данных - 1 часть

- Семейство MPP баз - назначение и особенности
- Логический и физический дизайн
- Vertica

3	<b>DWH. Хранилища данных - 2 часть</b>	<ul style="list-style-type: none"> <li>- Семейство MPP баз - назначение и особенности</li> <li>- Логический и физический дизайн</li> <li>- Google BigQuery</li> </ul>
		Домашние задания
		<p>1 Домашнее задание: проектирование витрины в Vertica (BigQuery).</p>
		<p>Цель: Спроектировать схему данных + Построить витрину          Использовать Vertica (Docker) или BigQuery • Датасет: Захват данных из divolte (или GCP Public Datasets) Definition of Done: • DDL объектов • DML шагов преобразовании • Опционально: Тестирование на наличие ошибок в данных</p>
4	<b>Хранилища NoSQL. Назначение и особенности.</b>	<ul style="list-style-type: none"> <li>- NoSQL Databases. HBase, Cassandra, Elasticsearch, Aerospike</li> <li>- Key-value</li> <li>- Cache</li> </ul>
5	<b>SQL-доступ к данным. Apache Hive.</b>	Домашние задания
		<p>1 HiveQL</p>
		<p>Цель: Практика с Hive на CDH</p>
		<p>Развернуть дистрибутив CDH          Самостоятельно проделать манипуляции с Hive с приложенными скриптами</p>
6	<b>Confluent Platform</b>	<ul style="list-style-type: none"> <li>- Apache Kafka &amp; Confluent platform</li> <li>- Schema registry. Данные с фиксированной схемой.</li> <li>- KStreams. Фреймворк для потоковой обработки.</li> <li>- KSQL. SQL на потоках данных.</li> </ul>
7	<b>Elasticsearch</b>	<ul style="list-style-type: none"> <li>- Знакомство с компонентами ELK-стэка</li> <li>- Классы задач, для которых подходит ELK</li> </ul>

## 4 Процессинг и доступ к данным

- |   |  |   |
|---|--|---|
| 1 | <b>Apache Spark - 1 часть</b>                            | <ul style="list-style-type: none"><li>- Spark - что это и зачем он нужен</li><li>- API - RDD, Dataset, Dataframe, операции над распределенными коллекциями</li><li>- Процесс вычисления в Spark - task, stage, оптимизатор запросов</li></ul> <hr/> |
| 2 | <b>Apache Spark - 2 часть</b>                            | <ul style="list-style-type: none"><li>- Spark - что это и зачем он нужен</li><li>- API - RDD, Dataset, Dataframe, операции над распределенными коллекциями</li><li>- Процесс вычисления в Spark - task, stage, оптимизатор запросов</li></ul> <hr/> |
| 3 | <b>Spark Streaming</b>                                   | <ul style="list-style-type: none"><li>- Micro-batch обработка данных</li><li>- Классический Spark Streaming</li><li>- Structured Streaming</li><li>- Continuous processing</li></ul> <hr/>  |
| 4 | <b>Доступ к данным, ноутбуки. Explore and visualize.</b> | <ul style="list-style-type: none"><li>- Инструменты интерактивной аналитики</li><li>- Google Cloud Datalab</li><li>- Jupyter - интеграция с Apache Spark</li></ul> <hr/>  |
| 5 | <b>Обучение моделей. ML.</b>                             | <p>Пример построения модели</p> <p>Домашние задания</p> <ol style="list-style-type: none"><li>1 Задание: обучаем собственную модель.</li></ol>  |

1	<b>Оркестрация</b>	<ul style="list-style-type: none"> <li>- Как организовать многоэтапные процессы обработки данных</li> <li>- Инструменты оркестрации - Oozie, Airflow</li> </ul>
2	<b>Интеграция, тестирование, развертывание. CI / CD. DevOps.</b>	<ul style="list-style-type: none"> <li>- Культура DevOps</li> <li>- Работа в команде</li> <li>- CI / CD</li> <li>- Auto tests</li> </ul>
3	<b>Мониторинг</b>	<ul style="list-style-type: none"> <li>- Инструменты мониторинга - Prometheus, Zabbix, Graphite, Grafana</li> <li>- Специфика мониторинга процессов обработки данных</li> </ul> <p>Домашние задания</p> <p>1      Задание: развернуть и настроить инструменты мониторинга. Проанализировать текущие показатели.</p>
4	<b>Data Quality. Контроль качества данных, мастер-данные, Troubleshooting.</b>	<ul style="list-style-type: none"> <li>- Data Quality and Consistency. Качество данных. MDM</li> <li>- Ошибки в коде, логике, виды, последствия, как найти и устранить корневую причину</li> <li>- Вопросы поддержки. Support</li> <li>- Network, integration, data quality, system faults, disk space, executor out of memory, grants, access rights, security</li> <li>- Метрики качества. Контроль качества. Data Fix - как исправлять найденные ошибки</li> <li>- MDM: управление мастер-данными</li> </ul>
5	<b>Case studies. Кейсы компаний.</b>	<p>Углубленные вопросы оптимизации. Фишки. Примеры, разбор</p> <p>Домашние задания</p> <p>1      Задание: разработать проверки качества данных для витрины. Внедрить их автоматическое выполнение.</p> <p>На предложенных примерах попытаться выполнить устранение ошибок (Data fix).</p>
6	<b>Бонус. Дальнейшее развитие Hard skills + Soft skills.</b>	<ul style="list-style-type: none"> <li>- Где искать ответы на вопросы. Ресурсы. Как быстро разбираться и решать проблемы.</li> <li>- Benchmarking - умеем сравнивать инструменты для решения конкретных задач</li> <li>- Как грамотно составить резюме (CV) + proof-read резюме участников курса</li> <li>- Как развиваться в плане Soft skills, Hard skills. Contribution.</li> </ul> <p>Домашние задания</p> <p>1      Задание: подготовить резюме (CV), отрецензировать резюме товарища. Завести аккаунт в LinkedIn.</p>



- |   |  |  |
|---|--|--|
| 1 | <b>Вводное занятие по проектной работе</b> | <p>Слушатели курса смогут определиться с темой проекта (выбрать из предложенного списка или привести задачу из деятельности своей компании), получить понимание какие ресурсы им необходимо использовать для работы.</p> <p>Домашние задания</p> <p>1 Проектная работа</p> <hr/> |
| 2 | <b>Консультация по проектной работе</b>    | <p>Слушатели курса получают комментарии относительно прогресса проектной работы, ответы на вопросы, рекомендации по реализации.</p> <hr/>  |
| 3 | <b>Защита проектной работы</b>             | <p>По окончании занятия слушатели курса получают разбор проектов, комментарии и оценку своей работы.</p>   |