**Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity.**

# Subjective opinion

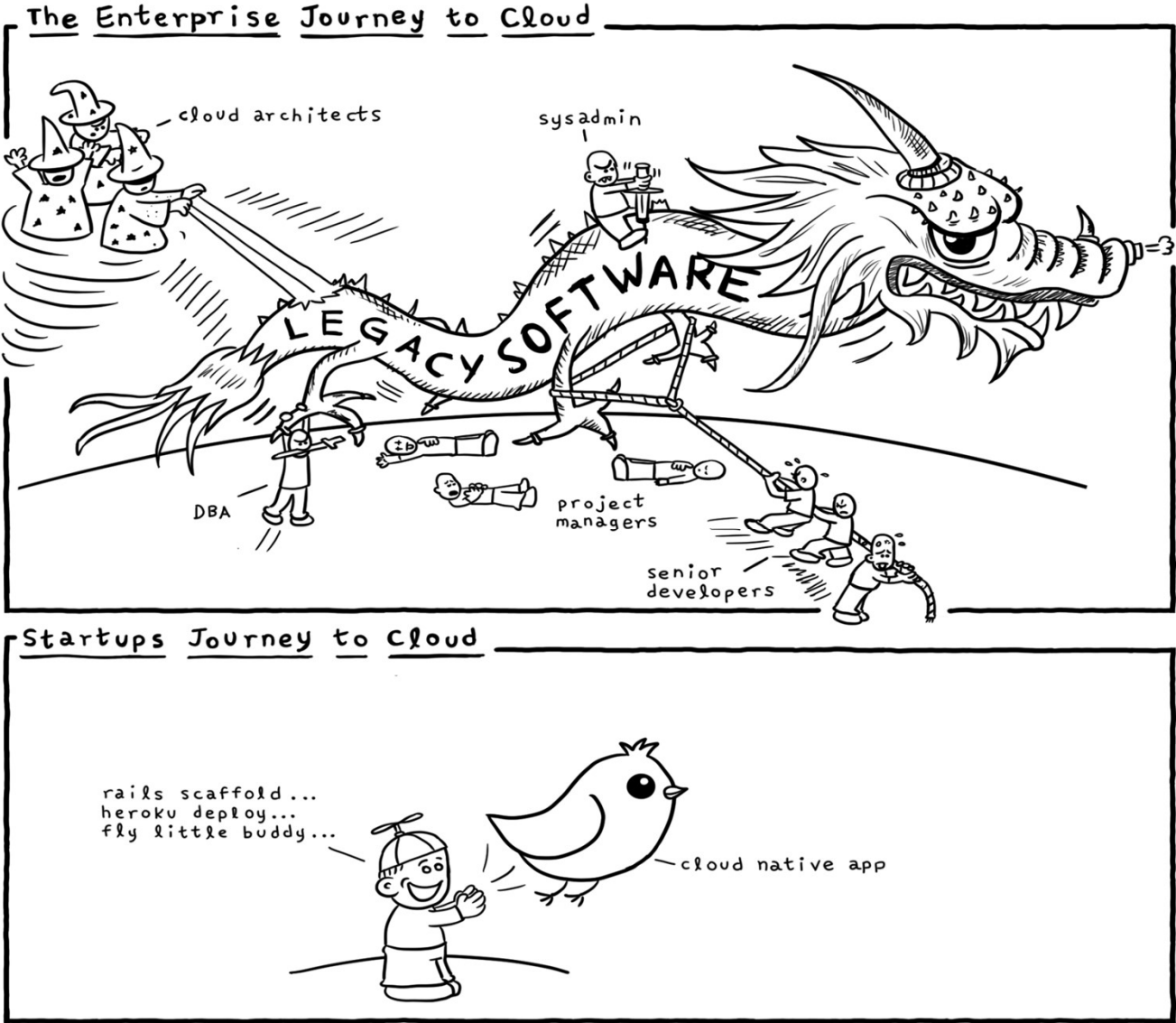Information from this report is my subjective opinion based on my experience, knowledge, mistakes... ;-)

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18        2

# Why you need emulate AWS

- Cost savings

- Testing automation. Isolated environments

- Improve development productivity. Network latency and throughput

- Great Firewall of China / Russia (AWS network blocking)

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18    3

# Marketing buzzwords is your enemy

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18    4

# Marketing buzzwords is your enemy

In economics, vendor lock-in, also known as proprietary lock-in or customer lock-in, makes a customer dependent on a vendor for products and services, unable to use another vendor without substantial switching costs.

https://en.wikipedia.org/wiki/Vendor_lock-in

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18          5

# Marketing buzzwords is your enemy

Amazon MQ is a managed message broker service that provides compatibility with many popular message brokers. **We recommend Amazon MQ for migrating applications from existing message brokers** that rely on compatibility with APIs such as JMS or protocols such as AMQP, MQTT, OpenWire, and STOMP.

Amazon SQS and Amazon SNS are queue and topic services that are highly scalable, simple to use, and don't require you to set up message brokers. **We recommend these services for new applications** that can benefit from nearly unlimited scalability and simple APIs.

https://docs.aws.amazon.com/amazon-mq/latest/developer-guide/welcome.html

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18    6

# Marketing buzzwords is your enemy

**Modernizing** Your Data Warehouse on AWS

**Easy Data Loading**

Load virtually any type of data into Amazon Redshift from a range of data sources including Amazon Simple Storage Service (Amazon S3), Amazon DynamoDB, Amazon EMR, and AWS Data Pipeline.

https://aws.amazon.com/big-data/featured-partner-data-warehouse-modernization/

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18        7

# Marketing buzzwords is your enemy

Data Types: SMALLINT, INTEGER, BIGINT, DECIMAL, REAL, DOUBLE PRECISION, BOOLEAN, CHAR, VARCHAR, DATE, TIMESTAMP, TIMESTAMPTZ

https://docs.aws.amazon.com/redshift/latest/dg/c_Supported_data_types.html

Amazon Redshift is based on PostgreSQL 8.0.2.

https://docs.aws.amazon.com/redshift/latest/dg/c_redshift-and-postgres-sql.html

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18    8

# Dependency inversion principle

You project code should depends on abstract/common API not on concrete cloud provider API.

- High-level modules should not depend on low-level modules. Both should depend on abstractions.

- Abstractions should not depend on details. Details should depend on abstractions.

https://en.wikipedia.org/wiki/Dependency_inversion_principle

Examples: Slf4J, Apache jclouds, micrometer.io, Apache Beam, Contexts and Dependency Injection (CDI)

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18          9

# Emulate deployment infrastructure

- Right choice.

    Examles: Testcontainers + LocalStack

https://github.com/localstack/localstack

https://github.com/testcontainers/testcontainers-java-module-localstack

- Pragmatic approach — in-JVM process libraries.

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro18     10

# Emulate Simple Cloud Storage Service (s3)

Heavy approach with real distributed object storage:

Ceph Object Gateway S3 API

https://github.com/ceph/s3-tests

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18          11

# Emulate Simple Cloud Storage Service (s3)

https://github.com/gaul/s3proxy

```
--
19   @Configuration
20   @ComponentScan("com.github.igorsuhorukov.postgresql")
21   public class Ctx {
22
23       private static final String S3_URL = "http://127.0.0.1:9988";
24
25       public AutoCloseable s3Proxy() throws Exception {
26           BlobStoreContext storeContext = ContextBuilder.newBuilder("transient").
27                   build(BlobStoreContext.class);
28
29           S3Proxy s3Proxy = S3Proxy.builder().blobStore(storeContext.getBlobStore()).
30                   awsAuthentication(AuthenticationType.AWS_V4, key, secret)
31                   .endpoint(URI.create(S3_URL)).build();
32           s3Proxy.start();
33           return s3Proxy::stop;
34       }
35
36       public AwsClientBuilder.EndpointConfiguration endpointConfiguration(){
37           return new AwsClientBuilder.EndpointConfiguration(S3_URL, signingRegion: "eu-central-1");
38       }
39
```

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18    12

# Emulate Simple Queue Service (SQS)

https://github.com/adamw/elasticmq

```java
12    @Configuration
13    public class SqsContext {
14
15        public static final String REGION = "elasticmq";
16        public static final String ENDPOINT = "http://localhost:9324";
17
18        public AutoCloseable sqsEmulator() {
19            SQSRestServer server = SQSRestServerBuilder.start();
20            return server::stopAndWait;
21        }
22
23        public AmazonSQS s3Client() {
24            return AmazonSQSClientBuilder.standard()
25                    .withCredentials(new AWSStaticCredentialsProvider(
26                        new BasicAWSCredentials(accessKey, secretKey)))
27                    .withEndpointConfiguration(
28                        new AwsClientBuilder.EndpointConfiguration(ENDPOINT, REGION))
29                    .build();
30        }
```

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18    13

# Emulate Simple Queue Service (SQS): JMS

spring-jms

com.amazonaws:amazon-sqs-java-messaging-lib

```java
private final JmsTemplate outboundQueue;

@JmsListener(destination = "${sqs.queue_name}")
@SneakyThrows
public void processMessage(Task task){
    var taskProcessor = getTaskProcessor(task);
    if(taskProcessor!=null){
        Status responseMessage = taskProcessor.processTask(task);
        outboundQueue.convertAndSend(outQueueName, responseMessage);
    } else {
        log.error("Task processor is not found. Task type {}",
                task.getClass().getName());

    }

}
```

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18     14

# Emulate Simple Queue Service (SQS): JMS

org.springframework.boot:spring-boot-starter-artemis

org.apache.activemq:artemis-jms-server

```java
@Import(ArtemisAutoConfiguration.class)
public class TestMessagingContext implements ArtemisConfigurationCustomizer {

    @Autowired
    @Qualifier("jmsConnectionFactory")
    ConnectionFactory testConnectionFactory;

    @Override
    public void customize(Configuration configuration) {
        configuration.setTransactionTimeout(TimeUnit.MINUTES.toMillis( duration: 2));
        Map<String, AddressSettings> addressesSettings =
                configuration.getAddressesSettings();
        AddressSettings addressSetting = new AddressSettings();
        addressSetting.setMaxDeliveryAttempts(0);
        addressesSettings.put("*", addressSetting);
    }


    @Bean
    public ConnectionFactory getConnectionFactory(Credentials credentials, Sqs sqs) {
        return testConnectionFactory;
    }
}
```

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18    15

# Emulate RDS PostgreSQL

http://www.h2database.com

https://github.com/yandex-qatools/postgresql-embedded

CDI wrapper: https://github.com/igor-suhorukov/postgresql-runner

Expose postgresql-embedded as CDI component

Manage Spring framework PG lifecycle as AutoCloseable…

Maven PostgreSql resolver

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18    16

# Emulate RDS PostgreSQL

```java
@Configuration
@ComponentScan("com.github.igorsuhorukov.postgresql")
public class MockPostgres {

    @Bean
    public String postgresVersion() { return "9.6.3-1"; }

    @Bean
    public String postgresDownloadPath() { return "org.postgresql:postgresql-server"; }

    @Bean
    public IDownloader downloader() {
        return new MavenDownloader();
    }

    @Bean
    public String postgresDatabaseStoragePath() { return ""; }

    @Bean
    public boolean initPostgresqlDatabase(@Autowired IPostgresqlService postgresqlService){
        String jdbcConnectionUrl = postgresqlService.getJdbcConnectionUrl();
        postgresqlService.importFromFile(new File( pathname: "init_script.sql"));
```

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18    17

# Emulate RDS PostgreSQL

com.github.springtestdbunit:spring-test-dbunit

```java
@RunWith(SpringJUnit4ClassRunner.class)
@ContextConfiguration(classes = {MockContext.class})
@TestExecutionListeners({ DirtiesContextTestExecutionListener.class,
        DependencyInjectionTestExecutionListener.class, DbUnitTestExecutionListener.class })
@DbUnitConfiguration(databaseConnection={"dbUnitDataSource"}, databaseOperationLookup = TruncateAndResetSerialKey.class)
@DirtiesContext(classMode = DirtiesContext.ClassMode.AFTER_CLASS)
public class DaoTest {
    @Autowired
    private AccountService accountService;
    @Rule
    public ExpectedException expectedException = ExpectedException.none();

    @Test
    @ExpectedDatabase(value = "/accounts/expected/user-save.xml", assertionMode = DatabaseAssertionMode.NON_STRICT_UNORDERED)
    @DatabaseTearDown(type = DatabaseOperation.CLEAN_INSERT, value = "/cleanup-tables-order.xml")
    public void testUserAccountSave() throws Exception {
        accountService.save(DaoTest.class.getResourceAsStream( name: "/user.json"));
    }
    @Test
    @ExpectedDatabase(value = "/accounts/expected/admin-save.xml", assertionMode = DatabaseAssertionMode.NON_STRICT_UNORDERED)
    @DatabaseTearDown(type = DatabaseOperation.CLEAN_INSERT, value = "/cleanup-tables-order.xml")
    public void testAdminAccountSave() throws Exception {
        accountService.save(DaoTest.class.getResourceAsStream( name: "/admin.json"));
    }
    @Test
    @DatabaseTearDown(type = DatabaseOperation.CLEAN_INSERT, value = "/cleanup-tables-order.xml")
    public void testDifCount() throws Exception {
        expectedException.expectMessage( substring: "phone should be provided in international format +(COUNTRY_CODE)NUMBER");
        expectedException.expect(IllegalArgumentException.class);
        accountService.save(DaoTest.class.getResourceAsStream( name: "/invalid-account.json"));
    }
```

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18    18

# Emulate Amazon Redshift (COPY)

https://github.com/opt-tech/redshift-fake-driver

This driver supports:

- json+jsonpath

- Manifest (json file with references to CSV files)

After source code modification driver is almost ready for CSV import

DriverClass: jp.ne.opt.redshiftfake.postgres.FakePostgresqlDriver

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro18    19

# Emulate Amazon Redshift (COPY)

```java
@Configuration
@ComponentScan("com.github.igorsuhorukov.postgresql")
public class RedshiftEmulator {

    @Bean public int postgresPort() { return 32111; }
    @Bean public String postgresVersion() { return "9.6.3-1"; }
    @Bean public String postgresDownloadPath() {
        return "org.postgresql:postgresql-server";
    }

    @Bean public IDownloader downloader() { return new MavenDownloader(); }

    @Bean public String postgresDatabaseStoragePath() { return ""; }

    @Autowired private IPostgresqlService postgresqlService;

    @Bean
    @SneakyThrows
    public DataSource redshiftDataSource(){
        String dbURL = postgresqlService.getJdbcConnectionUrl().
                replaceAll(regex: "postgresql", replacement: "postgresqlredshift");
        return ConnectionUtil.getDataSource(dbURL,
                postgresqlService.getUsername(), postgresqlService.getPassword(),
                dbDriver: "jp.ne.opt.redshiftfake.postgres.FakePostgresqlDriver", schema,
                CONNECTION_TIMEOUT, maxPoolSize: 5,
                isAutoCommit: true, registry: null, poolName: "redshift-dbcp",
                IDLE_TIMEOUT, VALIDATION_TIMEOUT);
    }
}
```
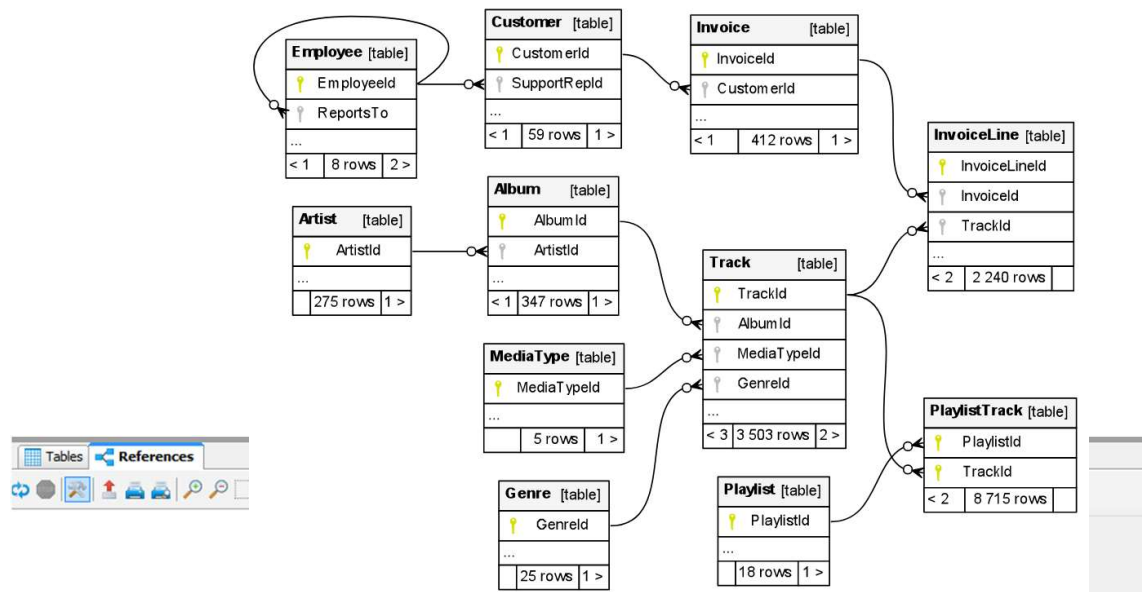
Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18    20

# Visualize Redshift table relations

https://www.dbvis.com/

http://schemaspy.org

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18    21

# Redshift JDBC «under the hood»
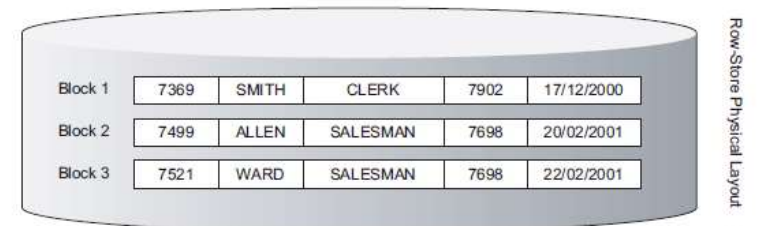


https://habrahabr.ru/post/345542/

Postgresql wire protocol 8.x

Implements jdbc:redshift and jdbc:postgresql

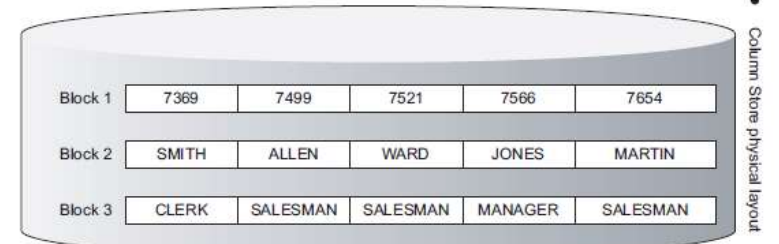Fat jar driver (packaged AWS SDK) is not worked with spring boot jar applications

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18        22

# Analytics database

## Column-oriented DBMS

| Operation | Row-oriented | Column-oriented |
|---|---|---|
| Aggregate operations | slow | fast |
| Insert/Update | fast | slow |
| Select single record | fast | slow |
| Select few columns | skip unnecessary data | fast |
| Compression | high | low |

## Data Lake



Row-Store Physical Layout

| Block 1 | 7369 | SMITH | CLERK | 7902 | 17/12/2000 |
| Block 2 | 7499 | ALLEN | SALESMAN | 7698 | 20/02/2001 |
| Block 3 | 7521 | WARD | SALESMAN | 7698 | 22/02/2001 |

Row Database stores row values together

Logical Schema

| EmpNo | EName | Job | Mgr | HireDate |
|---|---|---|---|---|
| 7369 | SMITH | CLERK | 7902 | 17/12/1980 |
| 7499 | ALLEN | SALESMAN | 7698 | 20/02/1981 |
| 7521 | WARD | SALESMAN | 7698 | 22/02/1981 |
| 7566 | JONES | MANAGER | 7839 | 2/04/1981 |
| 7654 | MARTIN | SALESMAN | 7698 | 28/09/1981 |
| 7698 | BLAKE | MANAGER | 7839 | 1/05/1981 |
| 7782 | CLARK | MANAGER | 7839 | 9/06/1981 |

Column Store physical layout

| Block 1 | 7369 | 7499 | 7521 | 7566 | 7654 |
| Block 2 | SMITH | ALLEN | WARD | JONES | MARTIN |
| Block 3 | CLERK | SALESMAN | SALESMAN | MANAGER | SALESMAN |

Column Database stores column values together

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro 18      23

# Analytics database: Redshift

Based on Postgresql 8.0.2 fork (ParAccel MPP)

v8.0.2 – 2005-04-07

+ AWS service integration, AWS hosted/managed

+ Regular SQL JOINs, support subqueries etc

- constraints, transactions

- Too small function set

- Scaling up downtime

- «slow» inserts/data import

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18     24

# Analytics database: Greenplum

Based on Postgresql 9.0

v9.0 – 2010-09-20

+ Open source PG fork

+ Support complex SQL queries

+ Rich functionality

- Scaling up and maintenance downtime

- Fork with backport of new features

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro18    25

# Analytics database: CitusDB

Based on Postgresql 10.0

v10.0 – 2017-10-05

+ Open source PG extension

+ Use latest PG versions and leverage it recent features

+ Distributed transactions

+ Rebalance shards without downtime

- Does not allow subqueries in the WHERE clause

- JOIN a local and a distributed table

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18    26

# Analytics database: Druid

+ Highly optimized for web metrics tasks

+ Very high ingesting rate

- Does not yet have full support for joins

- Limited SQL support

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro18    27

# Analytics database: ClickHouse

+ Highly optimized for web metrics tasks

- There is no global query plan for distributed query execution.

- Does not yet have full support for joins

- Limited SQL support

https://github.com/yandex/ClickHouse

https://ruhighload.com/doc/clickhouse/development/architecture.html

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18    28

## Analytics database: CrateDB

Elasticsearch based

+ Full text search, GIS functions

+ Presto SQL parser, PG wire protocol

+ Blob storage

- Constraints, transactions

- Query optimization

- Hash join only for 2 tables

https://habr.com/post/323742/

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18    29

# Analytics database: PrestoDB

+ connectors for external data format

+ dozen of functions:  window functions, geo etc

+ transaction support

- primitive table statistics

- query S3 data only with Hive

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18    30

# Analytics database: Apache Drill

+ schema free SQL

+ query S3/HDFS data directly

- ???

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18    31

# Analytics database: Dremio

\+ query data from S3, Redshift, Elasticsearch

\+ support Apache Arrow

\- small OSS community

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18    32

# Analytics database: Apache HAWQ

+ Interactively query Hadoop data, natively via HDFS

+ Cost-Based Optimizer

- Incubating status

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18        33

# Analytics database comparison

| Database | Based on | JOIN large tables | Data lake | Full text search, geo data |
|----------|----------|-------------------|-----------|----------------------------|
| Redshift | Postgres 8.0.2 | Yes | Redshift Spectrum | No |
| Greenplum | Postgres 9.0 | Yes | Postgres FDW | Yes |
| CitusDB | Postgres extension | Yes | Postgres FDW | Yes |
| Druid | | No | | |
| ClickHouse | | No | | |
| CrateDB | Elasticsearch, PrestoDB, Postgres wire protocol | 2 tables | | Yes |
| PrestoDB | | Yes | Yes | Geo functions only |
| Apache Drill | | Yes | Yes | |
| Dremio | | Yes | Yes | |
| Apache HAWQ | Greenplum | Yes | Yes | |

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

JavaCro'18    34

Igor Sukhorukov
Emulate Amazon Web Services infrastructure in single JMV process to reduce development cost and improve productivity

35

JavaCro'18

# Thank you!

igor.suhorukov@gmail.com
github.com/igor-suhorukov