

AWS Glue: The True Story of One Transfiguration

Кто такая Alice?



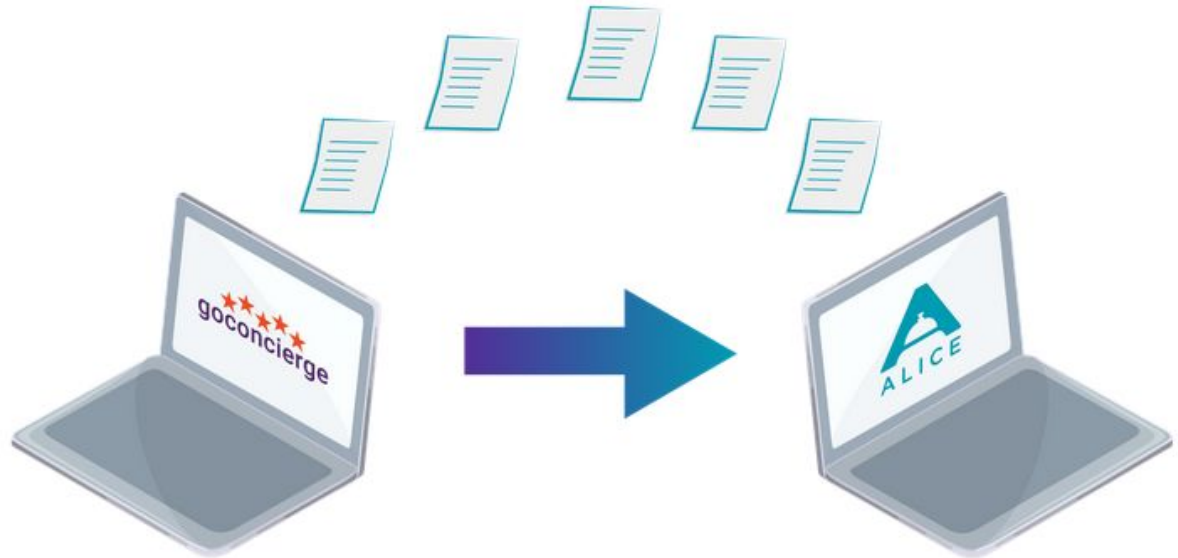
Before ALICE



After ALICE

GoConcierge

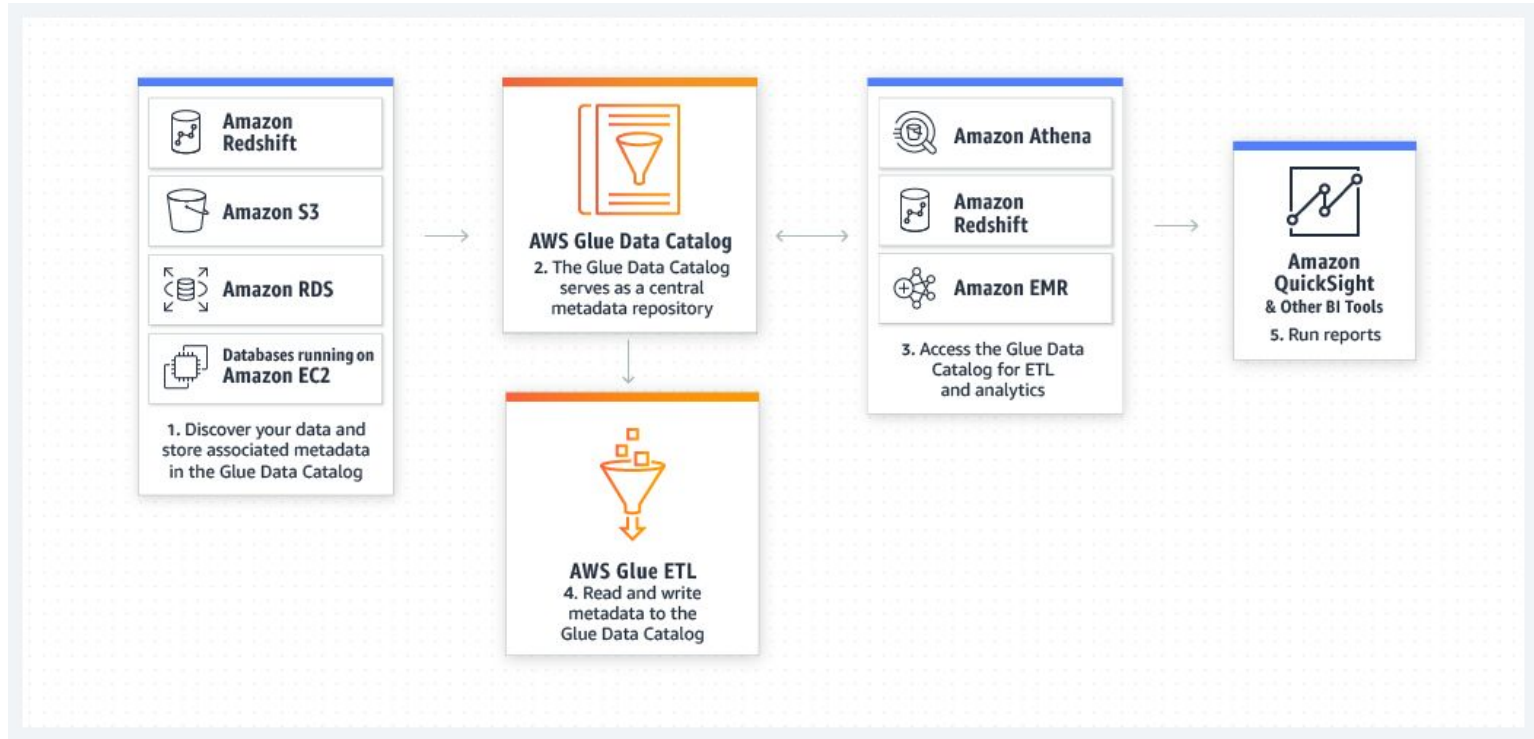
Sit back and relax! Upgrading to ALICE Concierge is a seamless installation. ALICE will import your GoConcierge users, tasks, task detail fields, vendors, and quick links.



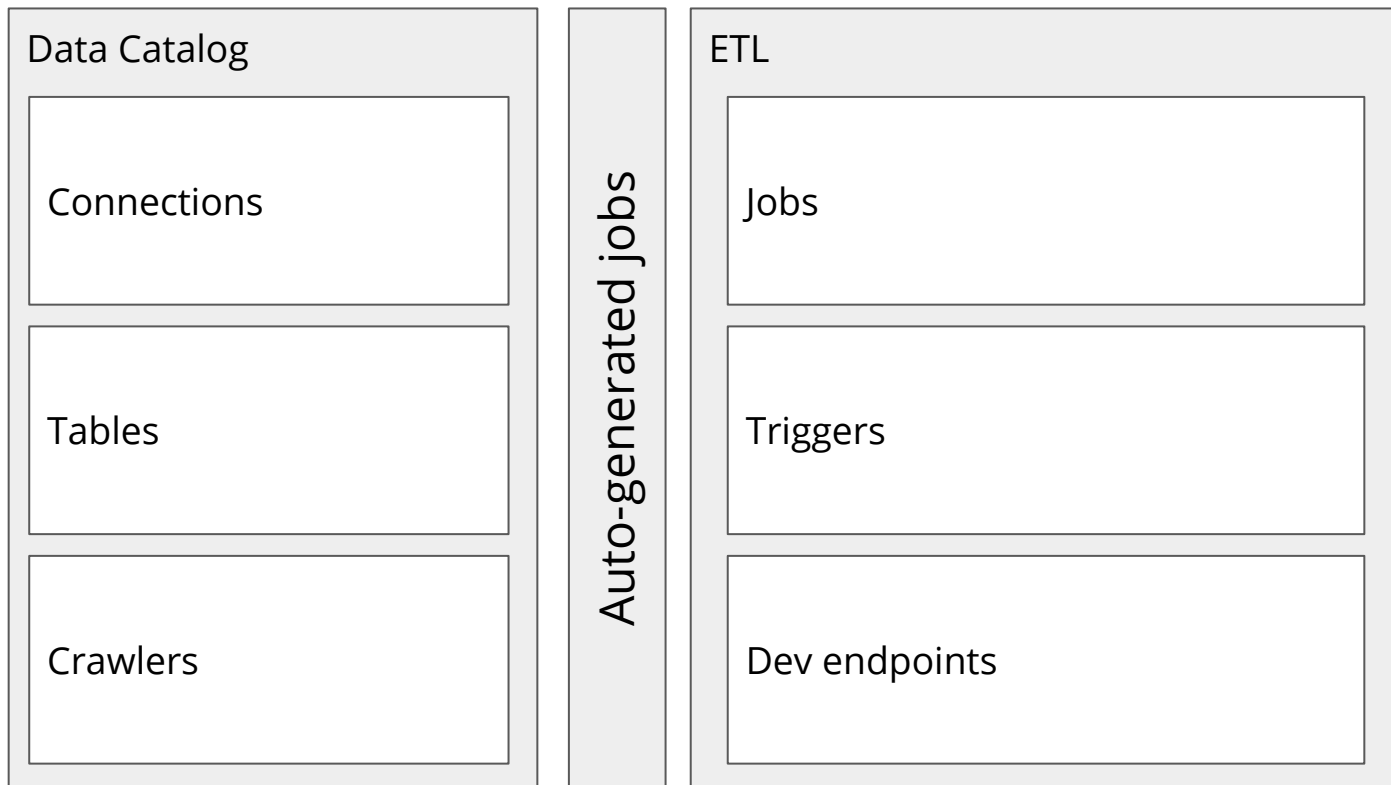
Выбираем ETL

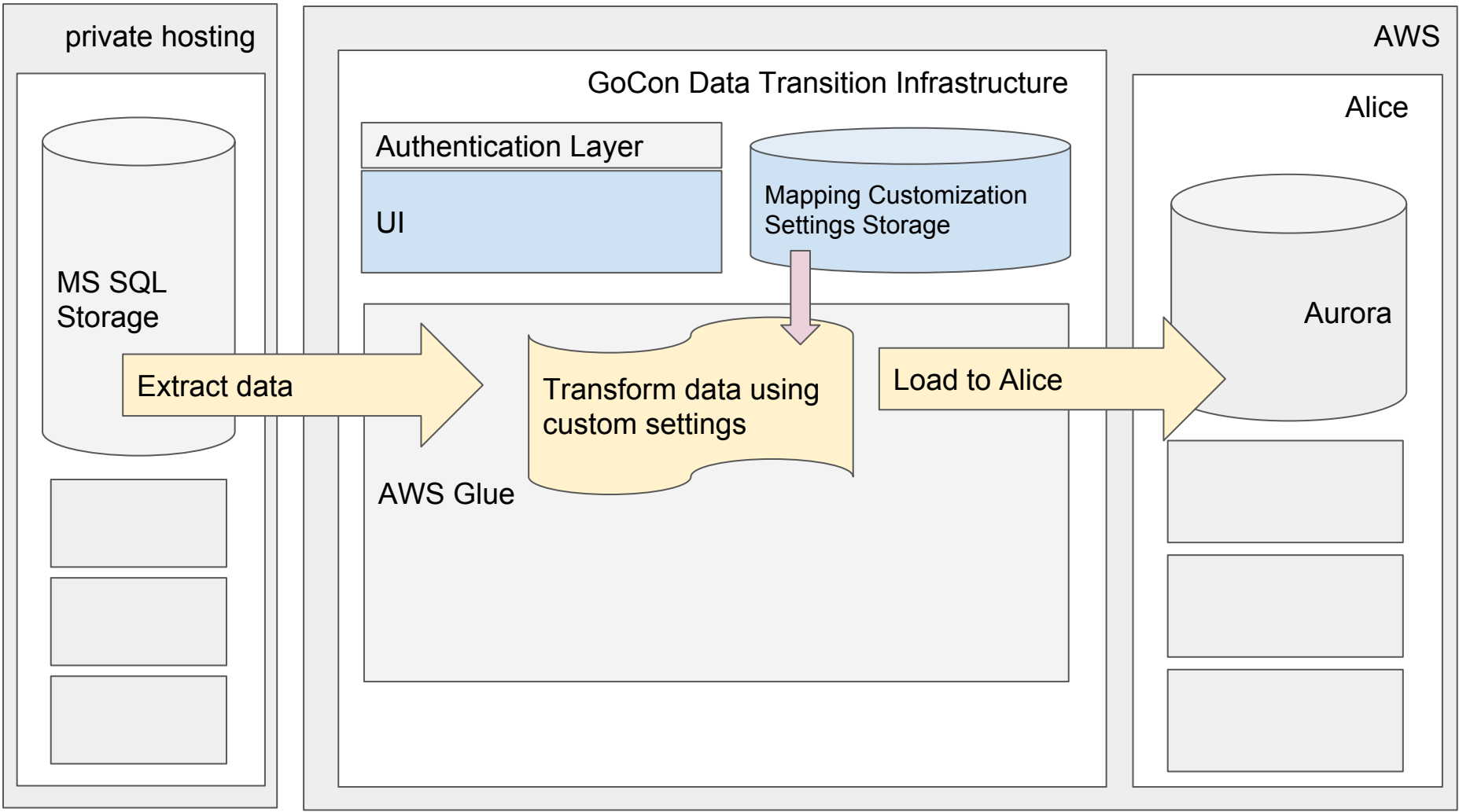


AWS Glue



Components





ROUND 1: Проще не придумаешь

Одна сущность - один Job

Сгенерируем скрипт, а чего не хватает - скопируем со Stackoverflow

Auto-generated jobs

The screenshot shows the AWS Glue console 'Add job' interface. The browser address bar displays `https://console.aws.amazon.com/glue/home?region=us-east-1#addJob:`. The page title is 'Add job'. On the left sidebar, there are navigation options: 'Job properties' (selected), 'Data source', 'Data target', 'Schema', and 'Review'. The main area is divided into 'Source' and 'Target' tables. The 'Source' table lists columns from the 'user' table, and the 'Target' table lists columns from the 'alice_user' table. Arrows indicate the mapping between the two tables.

Column name	Data type	Map to target
ext	string	-
loginname	string	lookup_key
bypasspw	boolean	-
firsttimepassword	boolean	-
department	string	-
accessitins	boolean	-
homehotel	int	-
edittask	boolean	-
otusermessage	string	-
accesstasksearch	boolean	-
changedbyid	int	-
bypassnotifypwd	boolean	password_expired
disablenotifications	boolean	-
accessreports	boolean	account_expired
emailbackup	boolean	-
emailaddress	string	email
tmp_encpassword_6	binary	-
allowrequestchange	boolean	-
ccpublic	boolean	enabled
pushtoweb	boolean	-
allowaddvendor	boolean	-

Column name	Data type
lookup_key	string
password_expired	boolean
account_expired	boolean
email	string
enabled	boolean
account_locked	boolean
version	long
username	string
id	long
password	string

© 2008 - 2018, Amazon Web Services, Inc. or its affiliates. All rights reserved. Privacy Policy Terms of Use

Auto-generated jobs

```
1 import sys
2 from awsglue.transforms import *
3 from awsglue.utils import getResolvedOptions
4 from pyspark.context import SparkContext
5 from awsglue.context import GlueContext
6 from awsglue.job import Job
7
8 ## @params: [JOB_NAME]
9 args = getResolvedOptions(sys.argv, ['JOB_NAME'])
10
11 sc = SparkContext()
12 glueContext = GlueContext(sc)
13 spark = glueContext.spark_session
14 job = Job(glueContext)
15 job.init(args['JOB_NAME'], args)
16
17 ## @type: DataSource
18 ## @args: [database = "insight_replica_glue", table_name = "insight_dbo_tbluser", transformation_ctx = "datasource0"]
19 ## @return: datasource0
20 ## @inputs: []
21 datasource0 = glueContext.create_dynamic_frame.from_catalog(database = "insight_replica_glue", table_name = "insight_dbo_tbluser", transformation_ctx = "datasource0")
22 ## @type: ApplyMapping
23 ## @args: [mapping = [{"loginname", "string", "lookup_key", "string"}, {"bypassnotifypwd", "boolean", "password_expired", "boolean"}, {"accessreports", "boolean", "account_expired", "boolean"}, {"emailaddress", "st
24 ## @return: applymapping1
25 ## @inputs: [frame = datasource0]
26 applymapping1 = ApplyMapping.apply(frame = datasource0, mappings = [{"loginname", "string", "lookup_key", "string"}, {"bypassnotifypwd", "boolean", "password_expired", "boolean"}, {"accessreports", "boolean", "acco
27 ## @type: SelectFields
28 ## @args: [paths = ["password", "account_expired", "account_locked", "password_expired", "lookup_key", "id", "version", "email", "enabled", "username"], transformation_ctx = "selectfields2"]
29 ## @return: selectfields2
30 ## @inputs: [frame = applymapping1]
31 selectfields2 = SelectFields.apply(frame = applymapping1, paths = ["password", "account_expired", "account_locked", "password_expired", "lookup_key", "id", "version", "email", "enabled", "username"], transformation
32 ## @type: ResolveChoice
33 ## @args: [choice = "MATCH_CATALOG", database = "alice_glue", table_name = "alice_user", transformation_ctx = "resolvechoice3"]
34 ## @return: resolvechoice3
35 ## @inputs: [frame = selectfields2]
36 resolvechoice3 = ResolveChoice.apply(frame = selectfields2, choice = "MATCH_CATALOG", database = "alice_glue", table_name = "alice_user", transformation_ctx = "resolvechoice3")
37 ## @type: ResolveChoice
38 ## @args: [choice = "make_cols", transformation_ctx = "resolvechoice4"]
39 ## @return: resolvechoice4
40 ## @inputs: [frame = resolvechoice3]
41 resolvechoice4 = ResolveChoice.apply(frame = resolvechoice3, choice = "make_cols", transformation_ctx = "resolvechoice4")
42 ## @type: DataSink
43 ## @args: [database = "alice_glue", table_name = "alice_user", transformation_ctx = "datasink5"]
44 ## @return: datasink5
45 ## @inputs: [frame = resolvechoice4]
46 datasink5 = glueContext.write_dynamic_frame.from_catalog(frame = resolvechoice4, database = "alice_glue", table_name = "alice_user", transformation_ctx = "datasink5")
47 job.commit()
```

```
alice_workflow_dyf = glueContext.create_dynamic_frame.from_catalog(database = 'alice_glue', table_name = 'alice_workflow')
```

```
alice_workflow_df = alice_workflow_dyf.toDF().select(['id', 'name']).withColumnRenamed('id', 'workflow_id')
```

```
alice_hotel_dyf = glueContext.create_dynamic_frame.from_catalog(database = 'alice_glue', table_name = 'alice_hotel')
```

```
alice_hotel_df = alice_hotel_dyf.toDF().select(['id', 'name']).withColumnRenamed('id', 'hotel_id')
```

```
alice_workflow_mapping_df = alice_workflow_df.join(alice_hotel_df, alice_workflow_df.name == alice_hotel_df.name)\
```

```
    .select(alice_workflow_df.workflow_id, alice_hotel_df.hotel_id)\
```

```
    .withColumn('version', lit(0))\
```

```
    .withColumn('deleted', lit(0))
```

```
alice_workflow_mapping_dyf = DynamicFrame.fromDF(alice_workflow_mapping_df, glueContext, "nested")
```

```
datasink = glueContext.write_dynamic_frame.from_catalog(frame = alice_workflow_mapping_dyf,
```

```
                database = "alice_glue",
```

```
                table_name = "alice_workflow_mapping",
```

```
                transformation_ctx = "datasink5")
```

```
job.commit()
```

Warning!

1. Долгий разогрев
2. Проблема связи между сущностями:

Mapping tables?

Нет однозначного соответствия между таблицами

- Одной GC таблице соответствует несколько в Alice и наоборот
- Joins мультиплицируют данные
- Сложно не запутаться

Lookup keys

EVENT_ID-APPOINTMENT_ID-WHAT_EVER_YOU_WANT_ID

```
# ----- hotel section -----
```

```
def add_hotel_lookup_key(self, dxf):
```

```
    return self.add_generic_lookup_key(dxf, ["CompanyID"])
```

```
def add_alice_hotel_reference(self, dxf):
```

```
    return self.add_generic_intermediate_reference(dxf, "alice_hotel", "hotel_id", ["CompanyID"])
```

ROUND 2: Сложнее, чем казалось

Python libraries

One job - many runs

Warning!

1. Очень долгие трансформации
2. Другие проблемы
 - a. Long warm up
 - b. Lazy filtering
 - c. Credentials
 - d. Mapping functions

ROUND 3: препроцессинг

JDBC + SQL

Оркестрация

CI/CD

Validation tests

```
join_df = spark.read \  
    .format("jdbc") \  
    .option("url", jdbc_url) \  
    .option("driver", "com.microsoft.sqlserver.jdbc.SQLServerDriver") \  
    .option("user", db_user) \  
    .option("password", db_password) \  
    .option("dbtable", query_string) \  
    .load()
```


?