# DESIGNING DATA-INTENSIVE APPLICATIONS IN SERVERLESS ARCHITECTURE

# Nikolay Matvienko

Senior Software Engineer and Node.js expert at Grid Dynamics

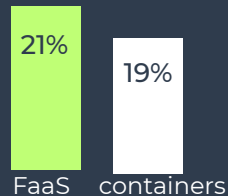Diagnostics, performance improvement consultant.

Work in USA and in Russia

Designed serverless computing platforms with AWS

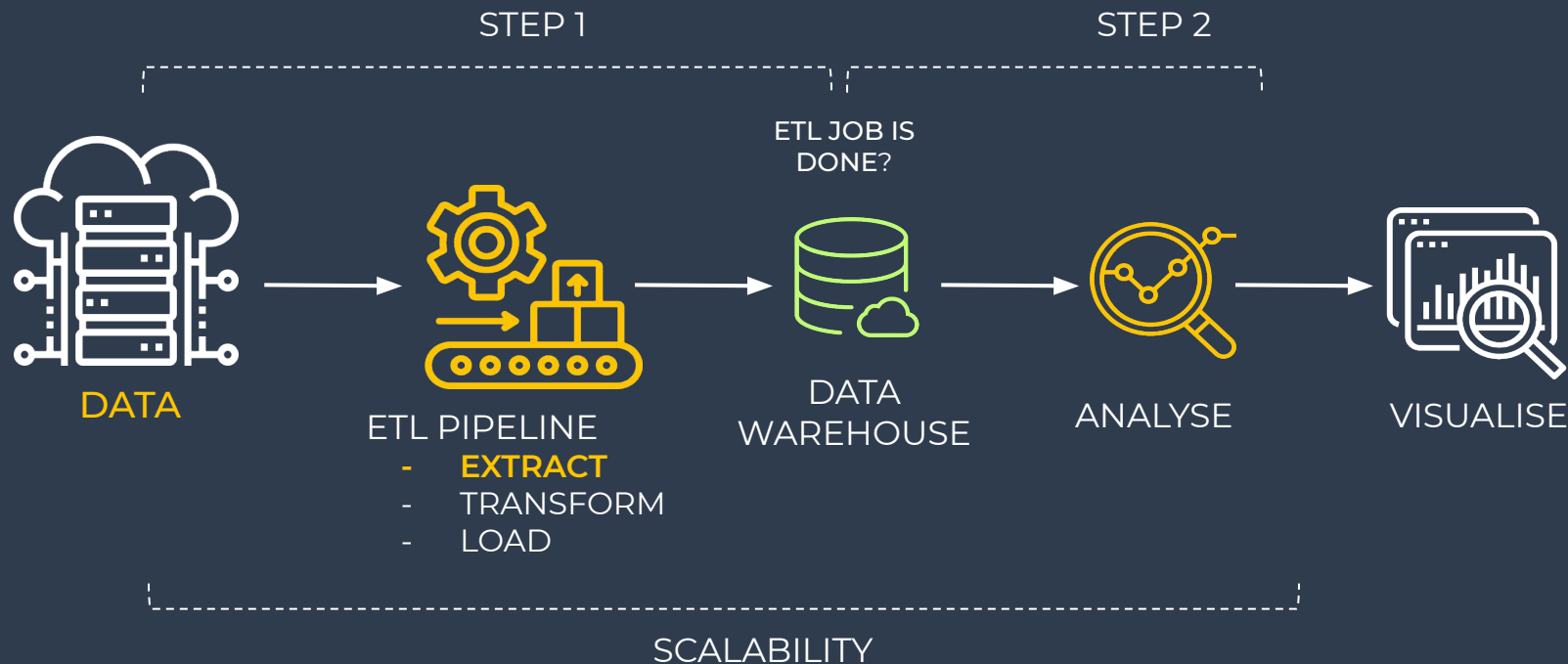You can find me at twitter.com/**matvi3nko** github.com/**matvi3nko**

2

# DESIGNING SERVERLESS

1. SERVERLESS COMPUTING HAS BECOME VERY POPULAR

2. SERVERLESS SOLUTIONS ARE ~60% AND UP TO SEVERAL TIMES CHEAPER

21%
FaaS

19%
containers

3. EVERY YEAR AWS ANNOUNCES NEW SERVICES

4. BUT THE LACK OF PATTERNS AND MISUSE CAN MAKE THE SOLUTION SEVERAL TIMES MORE EXPENSIVE

# DATA-INTENSIVE APPLICATIONS



STEP 1

STEP 2

DATA

ETL PIPELINE
- **EXTRACT**
- TRANSFORM
- LOAD

ETL JOB IS DONE?

DATA WAREHOUSE

ANALYSE

VISUALISE

SCALABILITY

@matvi3nko

4

# WHERE DOES THE DATA COME FROM?

**IoT** DEVICES

**REAL-TIME STREAM** PROCESSING

DB

ENTERPRISE PLATFORMS

**BATCH / STREAM** PROCESSING

API

API

DATA LAKE

DB

@matvi3nko

5

# USE EXISTING OR BUILD YOUR OWN PIPELINE

DATA

**ETL BATCH** PROCESSING    **STREAM** PROCESSING

GLUE JOBs

KINESIS STREAMS & ANALYTICS

AWS LAMBDA

DATA
WAREHOUSE

# DYNAMICALLY UPDATED PROCESSING LOGIC



PROCESSING RULES

PROCESSING RULES

DATA

WAREHOUSE

PROCESSING RULES

PROCESSING RULES

PROCESSING RULES

PROCESSING RULES

PROCESSING RULES

AWS LAMBDA

REAL-LIFE CASE:
**300+ CUSTOMERS**

EACH ESTABLISHES HIS
**OWN DATA
PROCESSING RULES**

@matvi3nko

7

# DYNAMICALLY UPDATED PROCESSING LOGIC

PROCESSING RULES

PROCESSING RULES

DATA

WAREHOUSE

PROCESSING RULES

PROCESSING RULES

PROCESSING RULES

PROCESSING RULES

PROCESSING RULES

AWS LAMBDA

REAL-LIFE CASE:
**300+ CUSTOMERS**

EACH ESTABLISHES HIS
**OWN DATA
PROCESSING RULES**

@matvi3nko

8

# SERVERLESS COMPUTE SERVICE

AWS **LAMBDA**

SERVERLESS FUNCTION

```javascript
export const handler = async (event) => {
  const data = event.Records[0].body;

  // - TRANSFORM data
  // - WRITE to DB or
  // - PUT TO QUEUE/STREAM
  return 'success';
};
```

- CHOOSE YOUR CODE LANGUAGE

- NO INFRASTRUCTURE TO MANAGE

- TRIGGERED BY EVENTS

- HIGH SCALABLE

- STATELESS

- COST-EFFECTIVE

\* Symbol in the presentation

@matvi3nko

# AWS S3 FOR DATA LAKE



@matvi3nko

# BATCH DATA PROCESSING

AWS LAMBDA ICON    REPLACED ON

SNAPSHOT:   ------------------------------------------------------------> DONE!
/yyyy=2019/mm=09/dd=27/

- 📁 folder–0
- 📁 folder–1
- 📁 folder–2
- 📁 folder–3
- 📁 folder–4
- 📁 folder–5
- 📁 folder–6
- 📁 folder–7
- 📁 folder–8

TABLES/FILES

MERGED RECORD
from tables/files

DB

FUNCTION    FUNCTION    FUNCTION

JSON

SCAN FOLDERS
GET ALL FILES

READ FILES
MERGE
RECORDS

PROCESS/ANALYZE
MODEL

# TRANSPORT SERVICES

### FAN-OUT

NOTIFICATION TOPIC

- PUB/SUB
- NOT WORRIED ABOUT DELIVERY

### KINESIS STREAMS

### REAL-TIME PROCESSING

DATA STREAMS

- BATCH 1–10000 RECORDS
- NOT AUTOSCALE
- 6 MB MESSAGE SIZE LIMIT
- EXACTLY ONCE PROCESSING
- HOURLY COST

### BATCH PROCESSING

QUEUE

- BATCH 1–10 RECORDS
- AUTOSCALE
- 256 KB MESSAGE SIZE LIMIT
- AT LEAST ONCE PROCESSING

### DYNAMODB STREAMS

### DISTRIBUTED TRANSACTIONS

DYNAMODB NoSQL DB

- BATCH 1–1000 RECORDS
- AUTOSCALE
- 6 MB MESSAGE SIZE LIMIT
- EXACTLY ONCE PROCESSING
- REACT ON EACH CHANGE

@matvi3nko

# BATCH PROCESSING ARCHITECTURE



SNAPSHOT:
/yyyy=2019/mm=05/dd=25/

TABLES/FILES
PATHS

MERGED RECORD
from tables/files

JSON

FUNCTION

FUNCTION

FUNCTION

DB

NOTIFICATION TOPIC

QUEUE

QUEUE

QUEUE

SCAN
GET FILES FOR
EACH FOLDER

READ FILES
MERGE RECORDS

PROCESS/ANALYZE
MODEL

# REAL-TIME PROCESSING ARCHITECTURE

SNAPSHOT:
/yyyy=2019/mm=05/dd=25/

TABLES/FILES
for each folder

MERGED RECORD
from tables/files

FUNCTION

FUNCTION

FUNCTION

JSON

DB

NOTIFICATION TOPIC

QUEUE

QUEUE

DATA STREAMS

SCAN
GET FILES FOR EACH
FOLDER

READ FILES
MERGE RECORDS

PROCESS/ANALYZE
MODEL

@matvi3nko

14

# DATA EXTRACTION PATTERNS

1. MOVE FROM BIG DATA TO A LARGE NUMBER OF MESSAGES

2. USE THE QUEUES FOR MESSAGES, AND DATA STREAMS TO TRANSFER MODELS / LARGE COLLECTION

3. BUT DO NOT RUSH TO USE STREAMS. CHOOSE TRANSPORT FOR YOUR NEEDS

# DATA TRANSFORMATION

STEP 1             STEP 2

ETL JOB IS
DONE?

DATA

ETL PIPELINE
- EXTRACT
- **TRANSFORM**
- LOAD

DATA
WAREHOUSE

ANALYSE

VISUALISE

SCALABILITY

@matvi3nko

16

# TRANSFORM SECTION



SNAPSHOT:
/yyyy=2019/mm=05/dd=25/

TABLES/FILES
for each factory

MERGED RECORD
from tables/files

JSON

FUNCTION

FUNCTION

FUNCTION

DB

NOTIFICATION TOPIC

QUEUE

QUEUE

QUEUE

SCAN
GET FILES FOR EACH
FACTORY

READ FILES
MERGE RECORDS

PROCESS/ANALYZE
MODEL

# SEPARATE FUNCTIONS BY RESPONSIBILITY



@matvi3nko

# BOUNDARY CONTEXT



CONTEXT 1

MERGER

FILTER

DATA STREAMS

CONTEXT 2

DATA STREAMS

WRITER

WRITER

S3

SPLITTER

DATA STREAMS

QUEUE

CONTEXT 3

KINESIS DATA ANALYTICS

DB

CONTEXT 4

PROCESS

DATA STREAMS

# QUERIES & CONNECTIONS PROBLEMS

SELECT

UPDATE

SQL DB

INSERT

TRANSFORM 2

TRANSFORM 1

TRANSFORM 3

PROBLEMS:

1. MANY LAMBDAS HAS TO QUERY SQL DB
2. A LOT OF CONNECTIONS
3. A LOT OF DEPENDENCIES

@matvi3nko

# DB QUERY LOGIC ENCAPSULATION



SQL DB

GET

UPDATE

**API**

PUT

TRANSFORM 2

TRANSFORM 1

TRANSFORM 3

BENEFITS:

1. QUERIES & LOGIC IS HIDDEN BEHIND THE API
2. LESS CONNECTIONS, CONTROLLED CONNECTIONS POOL
3. LESS DEPENDENCIES

# RELIABILITY: RETRY STRATEGY
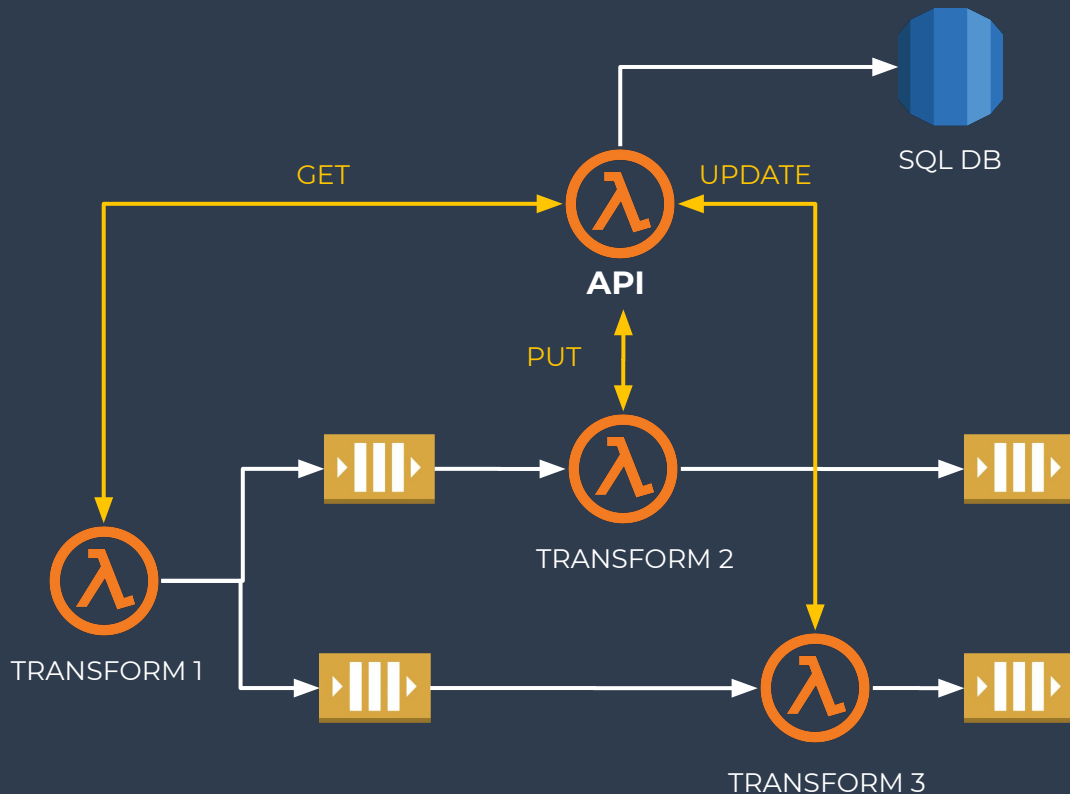


SCAN ALL OBJECTS FOR EACH FOLDER

GET ALL RECORDS AND AGGREGATE EACH WITH DETAILS

DATA LAKE

New SNAPSHOT DATE

QUEUE

DISPATCH

FILES PATHS

QUEUE

TRANSFORM

MODELS

QUEUE

WRITE

DATABASE

RETRIES 3 TIMES (BY DEFAULT)

logo

# DEAD LETTER QUEUE



@matvi3nko

23

# KINESIS ERROR HANDLING

RETRIES N TIMES
24h - 7d

ANALYZE
FUNCTION
...

INFRASTRUCTURE ERRORS

reject(error)

resolve()
SEND RECORD
TO SNS TOPIC

BUSINESS ERRORS
INTERNAL ERRORS

Pull by error type

REPROCESS
PULLS MESSAGES
FROM DLQ BY THE REQUEST

DLQ
STORES
FAILED RECORDS

@matvi3nko

# DLQ FOR THE QUEUE



SCAN ALL OBJECTS FOR EACH FOLDER

GET ALL RECORDS AND AGGREGATE EACH WITH DETAILS

DATA LAKE

New SNAPSHOT DATE

QUEUE

DISPATCH
MERGE BY FACTORY

OBJECTS PATHS

TOPIC

QUEUE

TRANSFORM

MODELS

QUEUE

ANALYZE

DATABASE

"reprocess_messages"

ACTOR

REPROCESSING TOPIC

REPROCESS
PULLS MESSAGES
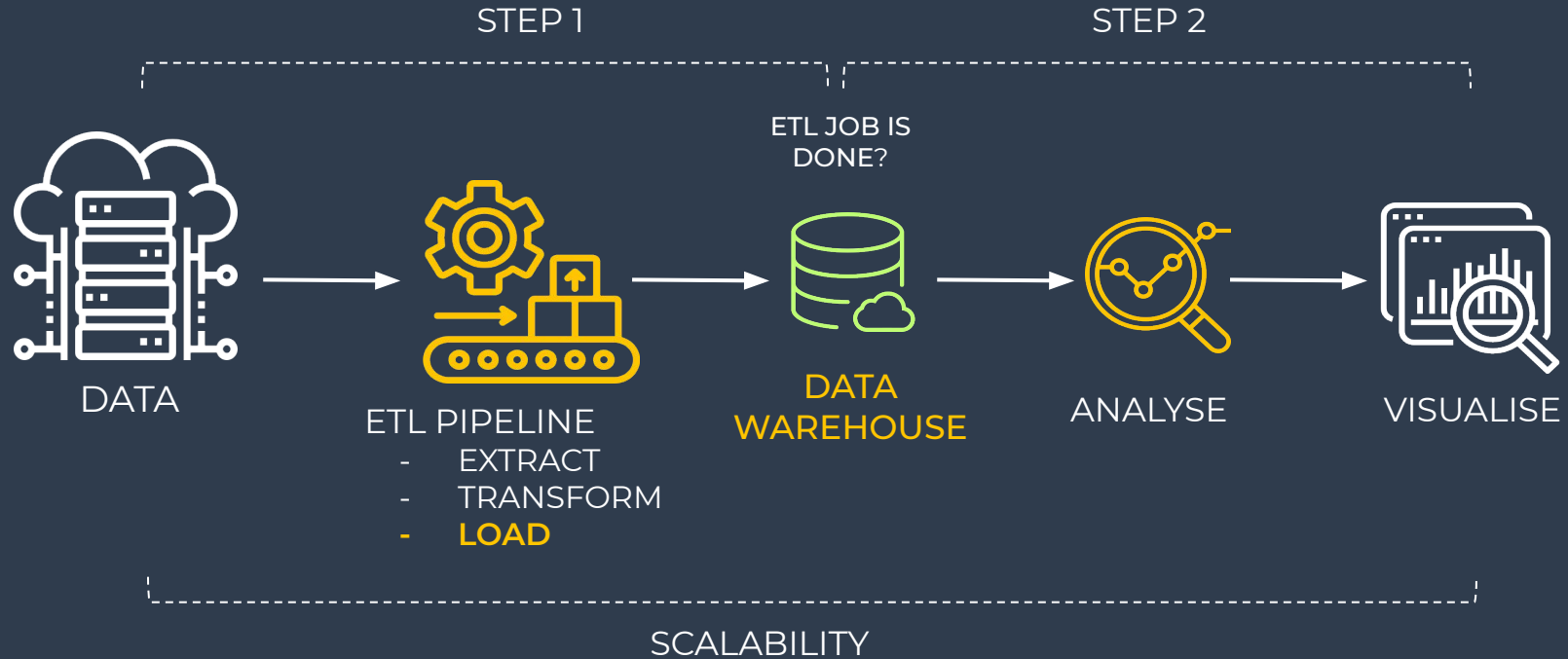FROM DLQ BY THE REQUEST

DLQ

1.  EASY TO TROUBLESHOOT

2.  NO NEED TO REPROCESS
    GIGABYTES OF DATA AGAIN

# DATA TRANSFORMATION PATTERNS

1. ONE LAMBDA FUNCTION – ONE RESPONSIBILITY

2. DIVIDE THE PIPELINE INTO BOUNDARY CONTEXTS WITH
   FIXED DATA INTERFACES

3. ENCAPSULATE DB QUERIES BEHIND API FUNCTION

4. USE DEAD LETTER QUEUES FOR RELIABILITY

@matvi3nko

# DATA LOADING



STEP 1                 STEP 2

ETL JOB IS DONE?

**DATA** — ETL PIPELINE — **DATA WAREHOUSE** — ANALYSE — VISUALISE

ETL PIPELINE
- EXTRACT
- TRANSFORM
- **LOAD**

SCALABILITY

# LOAD SECTION



SCAN ALL OBJECTS FOR EACH FOLDER

GET ALL RECORDS AND AGGREGATE EACH WITH DETAILS

DATA LAKE

New SNAPSHOT DATE

FILES PATHS

RECORDS

DATABASE

QUEUE

DISPATCH
MERGE BY TYPE

TOPIC

QUEUE

TRANSFORM

DATA STREAMS

ANALYZE & WRITE

"reprocess_messages"

DLQ

ACTOR

REPROCESSING TOPIC

REPROCESS
PULLS MESSAGES
FROM DLQ BY THE REQUEST

# WRITING TO A DATABASES



MICROSERVICE

Points of failure

1 – 10K
MESSAGES

LAMBDA

DB

⚡ PROCESSES DATA
ANR/OR FILTERS DATA

⚡ WRITES TO DB

Points of failure

1 – 10K
MESSAGES

LAMBDA

UI
DASHBOARD

ACCOUNTING

FURTHER DATA
PROCESSING

1. RETRY THE BATCH?
2. RETRY FAILED RECORD MANUALLY IN LAMBDA?

@matvi3nko

29

# RETRY PROBLEM



**1 – 10K**
MESSAGES

LAMBDA

UI
DASHBOARD

ACCOUNTING

FURTHER DATA
PROCESSING

DUPLICATE

DUPLICATE

1. [BAD] REMOVE/ROLLBACK SUCCESSFUL RECORDS

2. [BAD] ORCHESTRATE REPEAT IN LAMBDA

@matvi3nko

30

# DECOUPLED WRITERS



PROCESSING
ERRORS

PROCESSOR
/ FILTER

**1000**
INSTANCES

INFRASTRUCTURE
ERRORS

QUEUE    WRITER    DB 1

**200**
INSTANCES

INDEPENDENT
RETRIES

QUEUE    WRITER    DB 2

**100**
INSTANCES

BETTER SCALABILITY

QUEUE    WRITER    S3 DATA LAKE

**500**
INSTANCES

@matvi3nko

31

# DISTRIBUTED TRANSACTIONS



PROCESSOR / FILTER

QUEUE → WRITER → DB 1

QUEUE → WRITER → DB 2

QUEUE → WRITER → S3 DATA LAKE

HOW TO ROLLBACK OTHERS?

# USE AWS STEP FUNCTIONS



TRANSACTION

PROCESSOR / FILTER

AWS STEP FUNCTIONS

INSTRUCTIONS

WRITER — DB 1

WRITER — DB 2
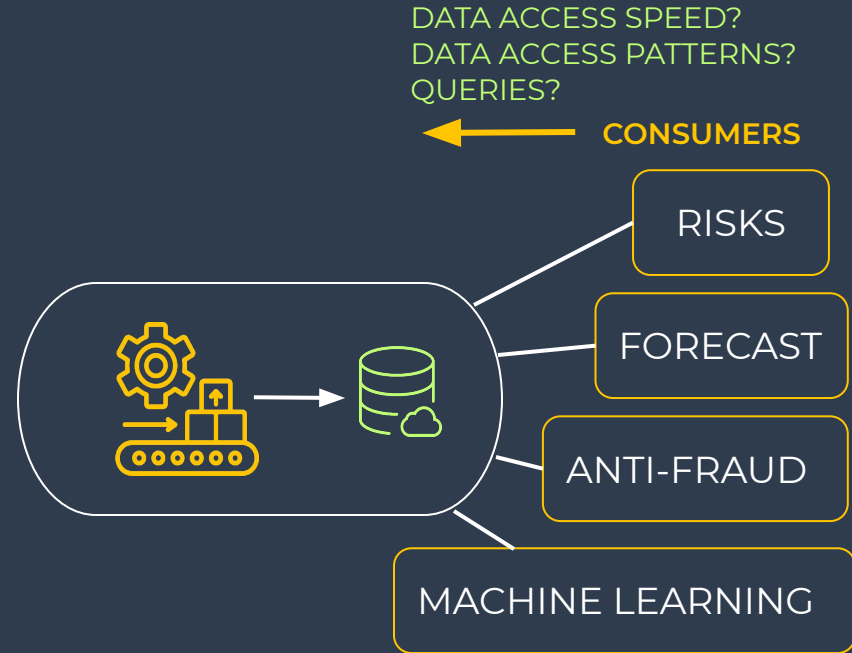
WRITER — S3 DATA LAKE

# DISTRIBUTED TRANSACTIONS

DYNAMODB PUSHES EVENT TO STREAM ON
INSERT | MODIFY | REMOVE
RECORD IN DB

WRITER 2 / STREAM SUBSCRIBER 1

S3 DATA LAKE

DYNAMODB STREAMS

PROCESSOR / FILTER

QUEUE

WRITER 1

DYNAMODB NoSQL DB

WRITER 3 / STREAM SUBSCRIBER 2

AURORA SQL DB

# DATA WAREHOUSE



S3
DATA LAKE

DYNAMODB
NoSQL DB

DB WRITER

AURORA
SQL DB

NEPTUNE
GRAPH DB

DATA ACCESS SPEED?
DATA ACCESS PATTERNS?
QUERIES?

CONSUMERS

RISKS

FORECAST

ANTI-FRAUD

MACHINE LEARNING

@matvi3nko

35
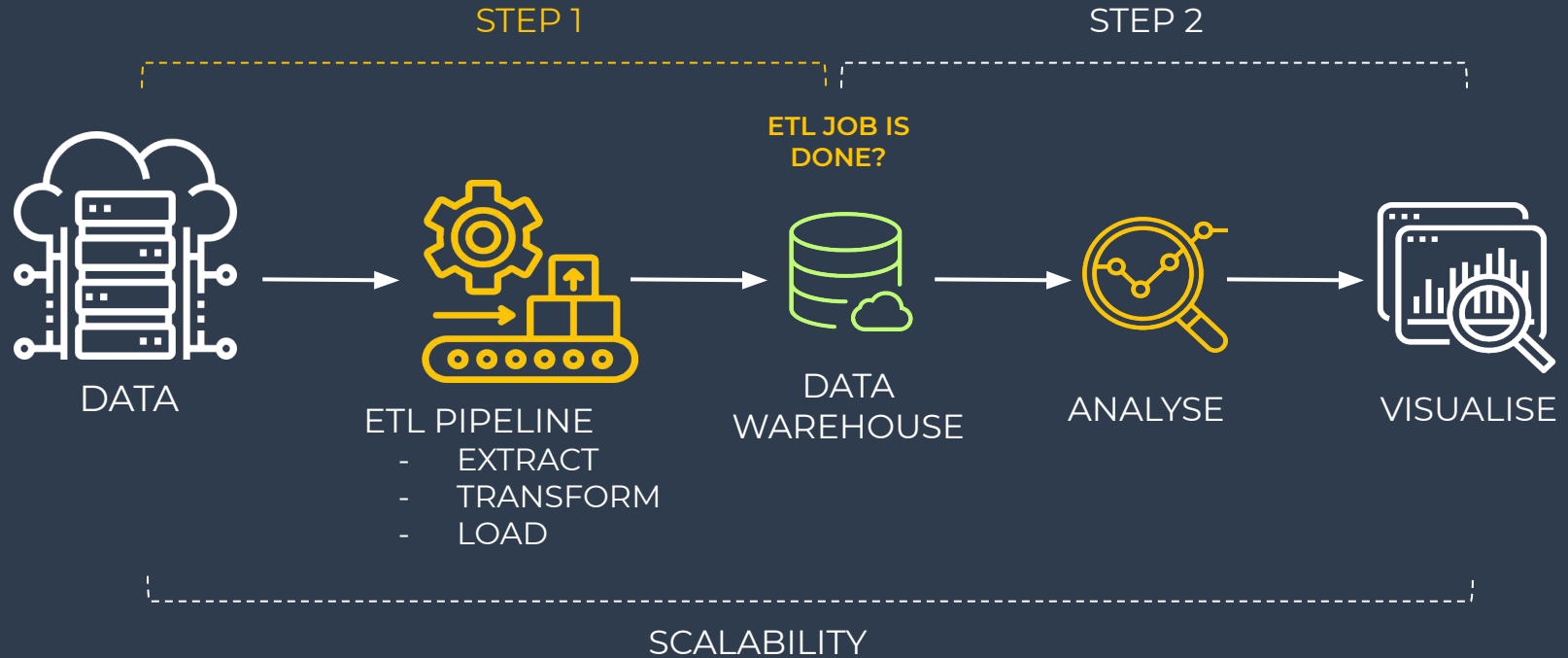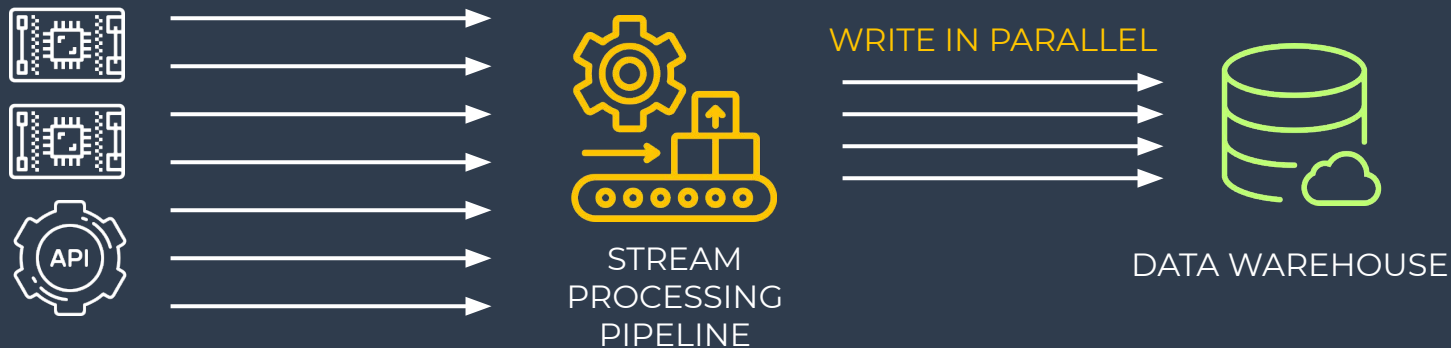
# DATA LOADING PATTERNS

1. DECOUPLE FUNCTIONS BY POINTS OF FAILURE

2. USE INFRASTRUCTURE AS A CODE VS ORCHESTRATIONS IN THE CODE

3. USE AWS STEP FUNCTIONS AND DYNAMODB STREAMS FOR TRANSACTIONS

4. PIPELINE IS THE CORE, THINK ABOUT FUTURE PLATFORMS AROUND IT AND

   HOW THEY WILL HAVE ACCESS TO DATA

@matvi3nko

36

# HOW TO UNDERSTAND THAT THE JOB IS COMPLETED?

STEP 1

STEP 2

ETL JOB IS DONE?

DATA

ETL PIPELINE
- EXTRACT
- TRANSFORM
- LOAD

DATA WAREHOUSE

ANALYSE

VISUALISE

SCALABILITY

@matvi3nko

37

# THE PROBLEM OF DISTRIBUTED DATA PROCESSING



DATA

ETL PIPELINE
EXTRACT TRANSFORM LOAD

WRITE IN PARALLEL

DATA WAREHOUSE

STREAM
PROCESSING
PIPELINE

WRITE IN PARALLEL

DATA WAREHOUSE

@matvi3nko

# ETL JOB STATE CONTROL

# DATA ANALYSIS



STEP 1

STEP 2

ETL JOB IS DONE?

DATA

ETL PIPELINE
- EXTRACT
- TRANSFORM
- LOAD

DATA WAREHOUSE

ANALYSE

VISUALISE

SCALABILITY
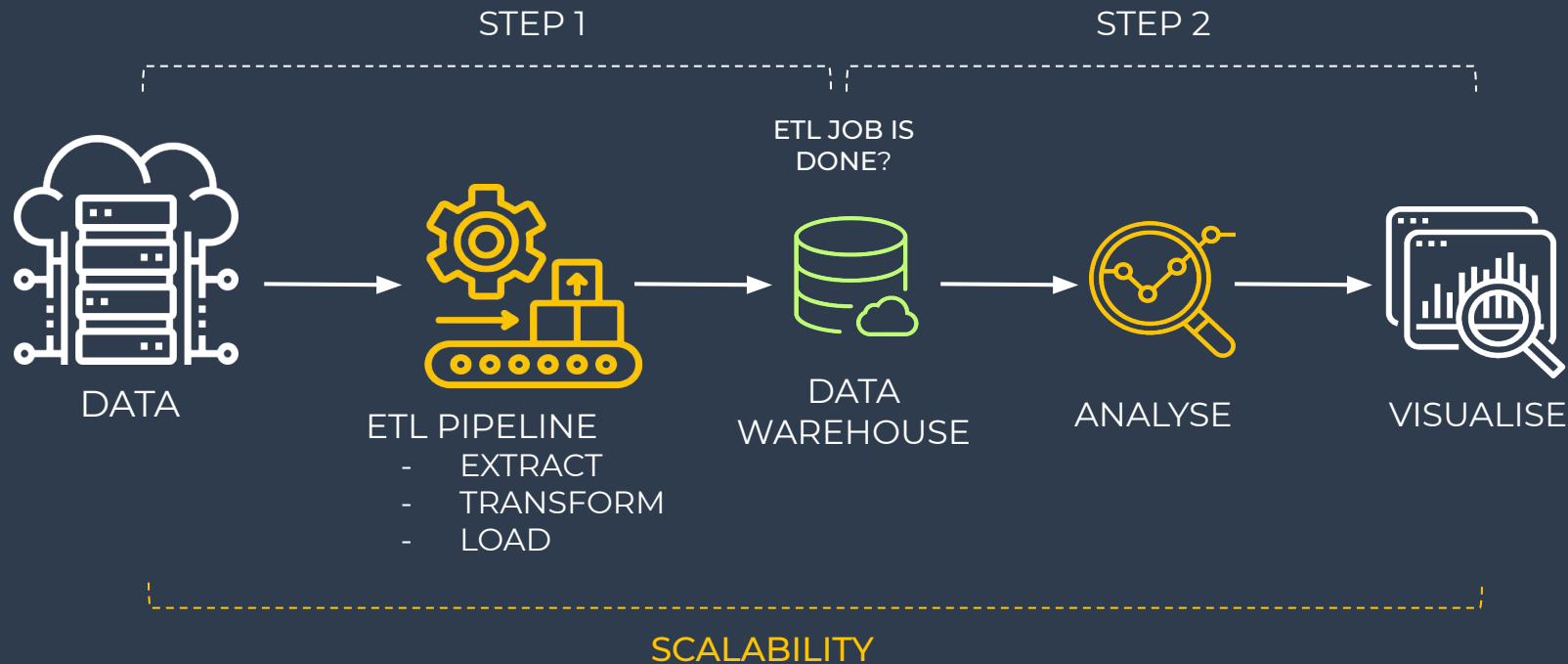
# ORCHESTRATE MULTIPLE ETL JOBS USING AWS STEP FUNCTIONS AND AWS LAMBDA



CLOUD WATCH
CHRON EVENT

SNS TOPIC
MESSAGE

AWS STEP
FUNCTIONS

INSTRUCTIONS

**STEP 1**

ETL JOB
RUNNER

GLUE JOB
ETL JOB

S3 DATA LAKE
STRUCTURED DATA

**STEP 2**

JOB STATE
CONTROLLER

ETL JOB STATE
DONE / NOT DONE

**STEP 3**

AWS LAMBDA ETL

READER

ANALYSER

@matvi3nko

41

# DATA ANALYSIS



STEP 1

STEP 2

ETL JOB IS
DONE?

DATA

ETL PIPELINE
- EXTRACT
- TRANSFORM
- LOAD

DATA
WAREHOUSE

ANALYSE

VISUALISE

SCALABILITY

# SCALABILITY



SNAPSHOT:
/yyyy=2019/mm=05/dd=25/

NOTIFICATION TOPIC    QUEUE    FUNCTION    QUEUE    FUNCTION    QUEUE    FUNCTION    DB

FILES

OBJECTS/FILES PATHS

DATE

~10

~100

~1000

ANALYZED REC

DB

@matvi3nko

43

# BLUEPRINT



ETL BATCH PROCESSING

ETL JOB RUNNER

GLUE JOB ETL JOB

S3 DATA LAKE STRUCTURED DATA

CLOUD WATCH CHRON EVENT

AWS STEP FUNCTIONS

JOB STATE CONTROLLER

ETL JOB STATE DONE / NOT DONE

SNS

JOB STATE

**5K – 10K FUNCTION** IN PARALLEL

PROCESSES **10-100 GB/SEC**

AD-HOC STREAM PROCESSING

READER

DYNAMO JOBS

DLQ

DLQ

PROCESSORS

DYNAMO

REDIS

KINESIS

AURORA

DLQ

44

@matvi3nko

# CONCLUSIONS

1. DESIGN MISTAKES = SOLUTION PRICE

2. IN SERVERLESS YOU CAN BUILD VERY BUILD A VERY FLEXIBLE ARCHITECTURE

   YOU CAN SWITCH FROM BATCH TO STREAM PROCESSING

3. YOU CAN RUN THOUSANDS OF LAMBDA FUNCTIONS IN PARALLEL

   AND PROCESS GBs OF DATA PER SEC

# THANKS!

## Nikolay Matvienko

Grid Dynamics
We are now in Belgrade too.

You can find me at twitter.com/**matvi3nko** github.com/**matvi3nko**