# ETL Quick Dive

A Guide by Ivan Sokolov with 💖

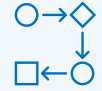# CONTENTS
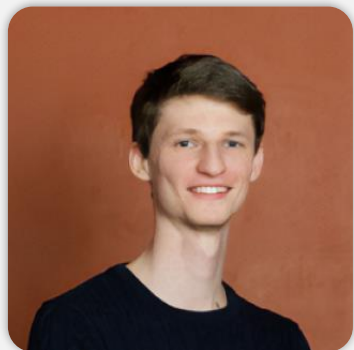
Hey there!
I'm Ivan.
I like automation, fixing data discrepancies, and piano music

# Tbilisi
Current location

## 3,5 years
Overall experience

## SQL
Main skill

# FMCG
# INSURTECH
Domain knowledge

**Atria®**
PERHETILOILTA VUODESTA 1903

→

**Symfa™**

**2019-2021**
Forecasting
Business Intelligence (BI)
ETL

SQL (SSMS)
IBM SPSS Modeler
Python

**2021-now**
BI – Enterprise Reporting
ETL – International DWH

SQL (SSMS, SSRS, SSIS)
Astera Centerprise
Git

ETL stands for 'Extract, Transform, Load'

This technology is very popular among large enterprises
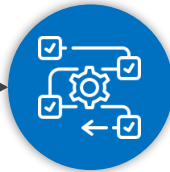
If you collect data from different counteragents, you need ETL

## Extract



DataBase

File / Document

ERP

## Transform



Cleaning / Grooming

Validation

Aggregation

Merging / Appending

Normalization

## Load



Data Warehouse (DWH)

Folder

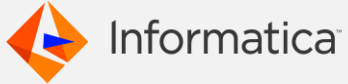Portal / Website

Cloud

## BI



## ML



Data Engineers, ETL Developers

BI Devs, Data Analysts, ML Devs

The most popular ETL tools
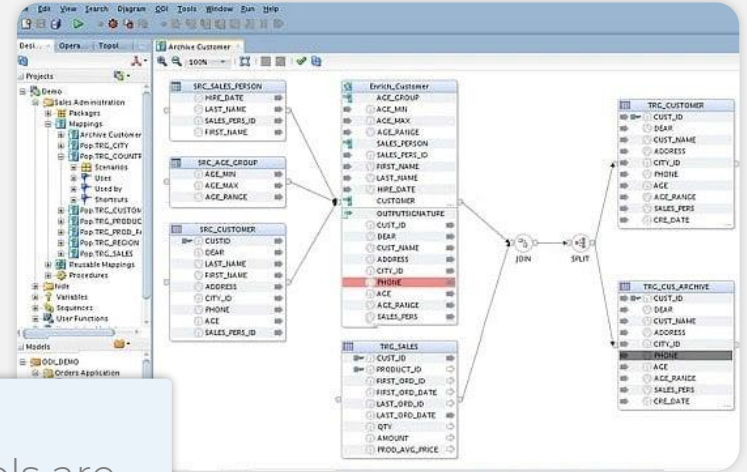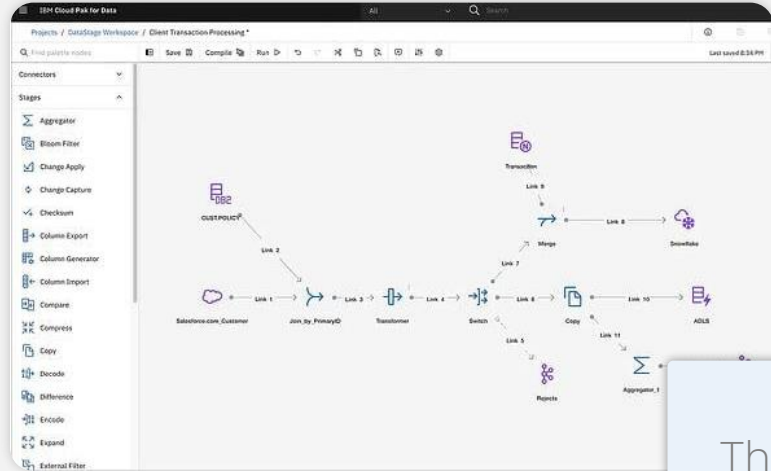
The major part of ETL tools are designed with similar UI

You know, I'm something of a scientist myself.

If you have ever merged / appended several Excel spreadsheets / built pivot tables and saved it to a new file, you already have some ETL experience



Spreadsheet

VLOOKUP

Spreadsheet

New Excel File

5

A **good** ETL Dev is a guy who:
- Understands and sees his data on the **granular** level
- Knows how to **manipulate** the data
- Understands what is going on behind **every step** of his flows

Create flows

Automate flows

Schedule flows

Update flows

## Required Skillset

**Tools / Software**
Microsoft, Amazon, IBM, Apache or any other tool

**SQL**
Querying the data

**Parametrization**
Dynamically change ETL aspects

**Scripting Language**
Juggle files, directories, users, and permission issues

**Organization**
Keep the work at hand organized and structured
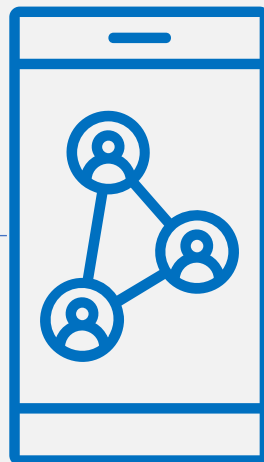
**Creativity**
Source To Target Mapping

**Debugging**
Fix the broken

UX / UI

Functionality

Interface

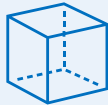Application

API

Server / DWH / Data

System Components

*ETL is here!*

ETL Development is **different** from other development types and has its **own features**

## Functional Development

**Software, app, web**
object of development

**CI/CD, Containers**
main tool / approach

**PL**
Ruby, JS, Python, Swift, C++, C, C#, Perl, R, SQL

**Coding Paradigm**
OOP, FP

## ETL Development

**Data**
object of development

**ETL Pipeline**
main tool / approach

**PL**
**SQL**, Python, R

**SQL Paradigm**
Declarative language
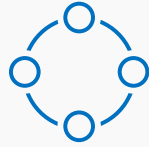
## Spot Check

Simple and fast queries to check the data at random places

```
SELECT TOP <N>
<…>
FROM <Table>
WHERE <...>
```

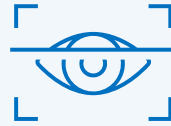*Quick but not reliable*

## Integrity Check

Heavy queries to check the relations between the tables of a database
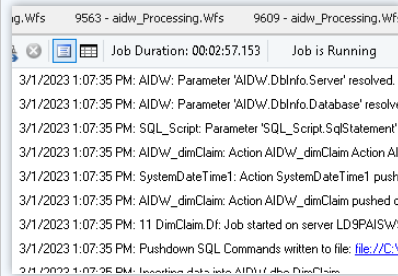
```
SELECT *
FROM <DimTable>
WHERE <Key> = -1


SELECT *
FROM <FactTable>
WHERE <Key> = -1
```

*Reliable but long*

## Empirical Test

Empirical analysis of the way the ETL works



*Semi-reliable*

## Execution Plan

SSMS feature that graphically represents operations' execution. Helps identify bottlenecks



*Powerful but needs rights access*

## Source-Destination Check (SDC)

SQL queries to check destination tables against their source by a particular criteria

```
SELECT
,SUM(<Field>)
,COUNT(*)
FROM <Destination>


SELECT
,SUM(<Field>)
,COUNT(*)
FROM <Source>
```

*Powerful and my favorite*

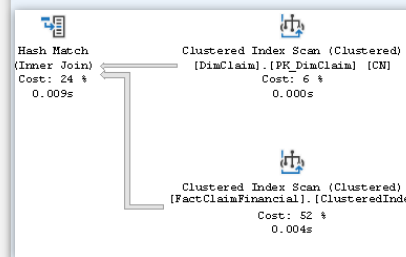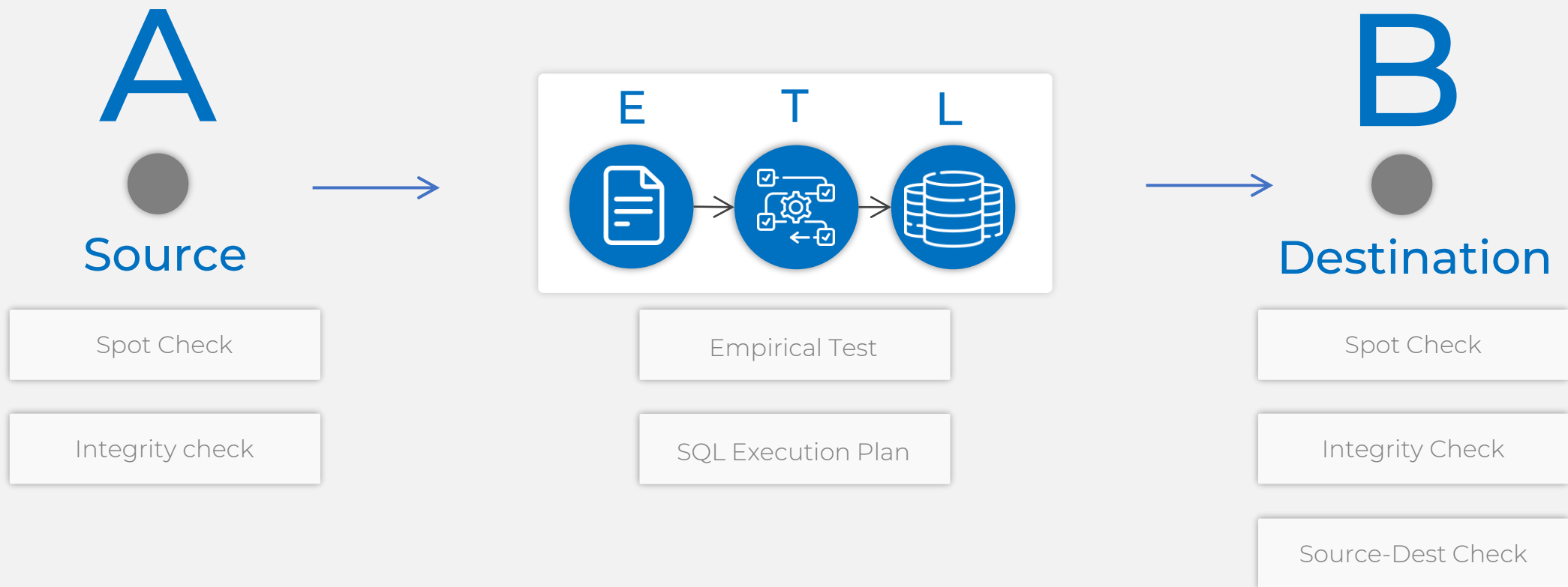Testing methods in a general dataflow

# A
### Source

# B
### Destination

E T L

Spot Check

Integrity check

Empirical Test

SQL Execution Plan

Spot Check

Integrity Check

Source-Dest Check

The workflow of applying testing methods

**1**

**ETL Test**

Empirical Test

↓

SQL Execution Plan

**2**

**Destination Test**

Integrity Check

↓

Spot Check

Source-Dest
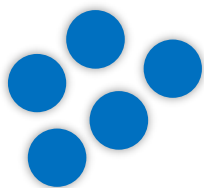
↓

Spot Check

**3**

**Source Test**

Spot Check

ETL Fix

Source Fix

11

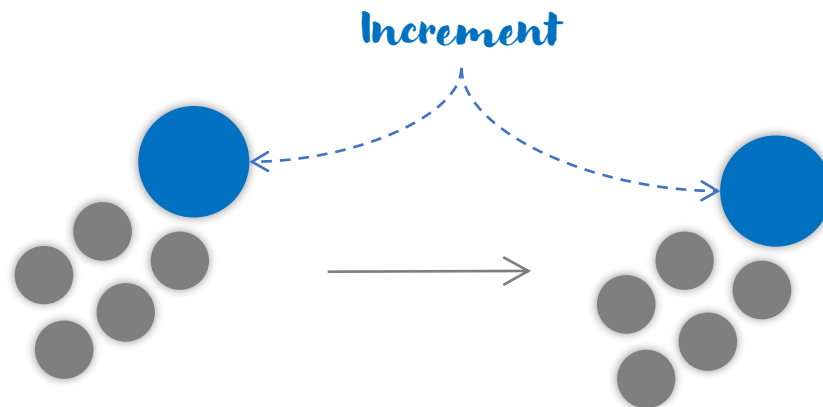When dealing with **large historical** data, at some point every ETL should use **incremental** loads



## Time Period 1

Source                    Destination

## Time Period 2

*Increment*

Source                    Destination

The closer you get to **Production**, the more **features** should be implemented and tested

| | | Scheduling | Scheduling | |
| | Logging | Logging | Logging | |
| | Full Automation | Full Automation | Full Automation | |
| | Auto Delete | Auto Delete | Auto Delete | |
| Incremental Load | Incremental Load | Incremental Load | Incremental Load | |
| ETL | ETL | ETL | ETL | |

DEV                 QA                 UAT                 Pre-PROD                 PROD