

# Инжиниринг Данных в Microsoft

*основано на реальных событиях*



SmartData

2022

# Disclaimer

All thoughts are mine and don't represent Microsoft or any other company.



# Содержание

- Про меня
- Собеседование и начало работы;
- Создание Delta Lake с нуля и модернизация аналитического решения;
- специфика Microsoft;
- плюсы и минусы работы в Microsoft;
- отличие культур Amazon и Microsoft;
- типовые аналитические решения Microsoft.

- 12 + лет в аналитики, до этого был инженером конструктором на ГКНПЦ им Хруничева
- Москва -> Черногория -> Канада
- 5 лет @Amazon, 2 года @Microsoft Gaming
- Tableau, Snowflake, Microsoft, AWS пользовательские группы и митапы, конференции



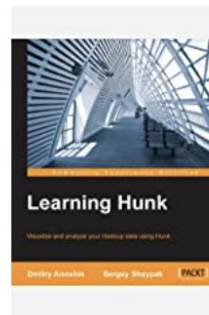
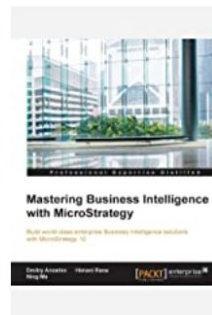
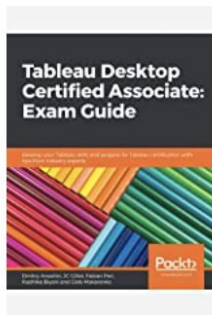
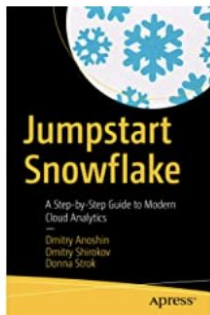
DataLearn.ru  
7000 Students

- DE 101
- DS&ML 101
- SQL 101



Инжиниринг Данных  
10 781 subscribers

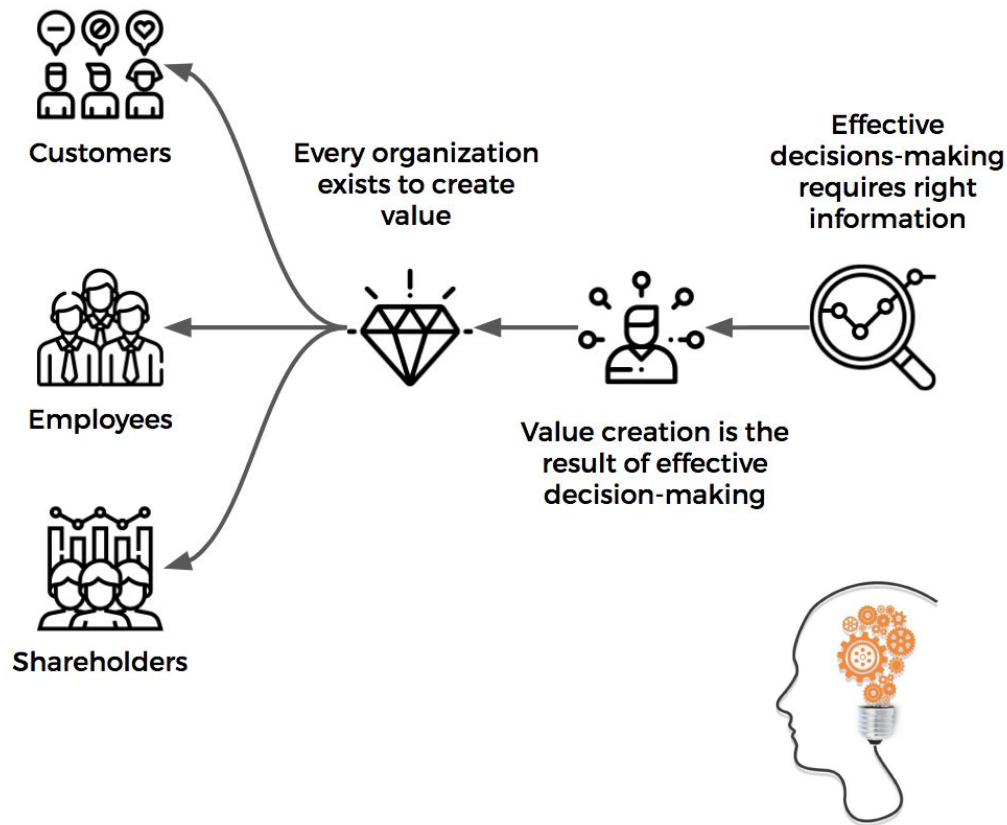
Инжиниринг данных  
<https://t.me/rockyourdata>



# Часть 1: Немного теории

# Зачем нужна аналитика?

- Повысить прибыль
- Сократить расходы
- Избежать риски
- Исследование новых рынков и продуктов
- Проверить гипотезы



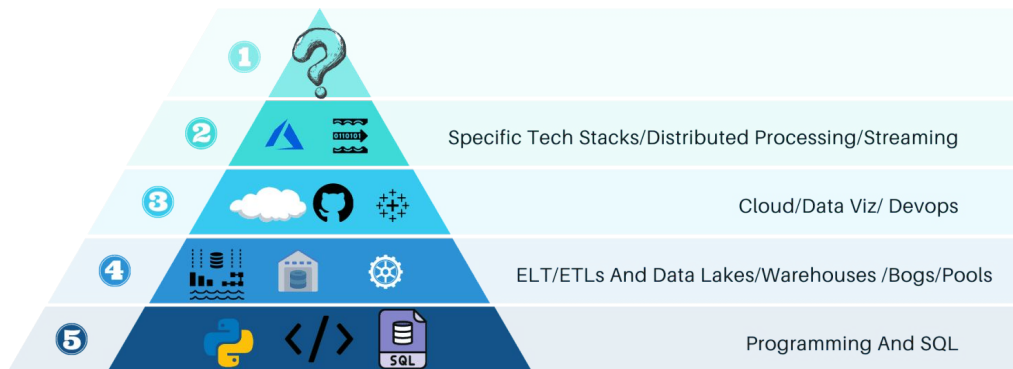
## HYPER

Changing the way you think about, plan, and execute  
**Business Intelligence**  
for real results, real fast!

Gregory P. Steffine

# Что такое Инжиниринг данных?

Data Engineering Hierarchy of Skill Sets



Инжиниринг данных делает данные **полезными** и **доступными** для потребителей, создавая **безопасную** и **масштабируемую** инфраструктуру данных.

# Modern Data Stack

## Source Layer



Files, SFTP,  
etc.



IoT



APIs



OLTP

## Data Processing



Batch  
(ETL/ELT)



Streaming

## Storage



Data Warehouse

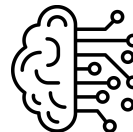


Big Data Solution



Data Lake

## Science & Experimentation



Datascience  
Machine  
Learning

## Business

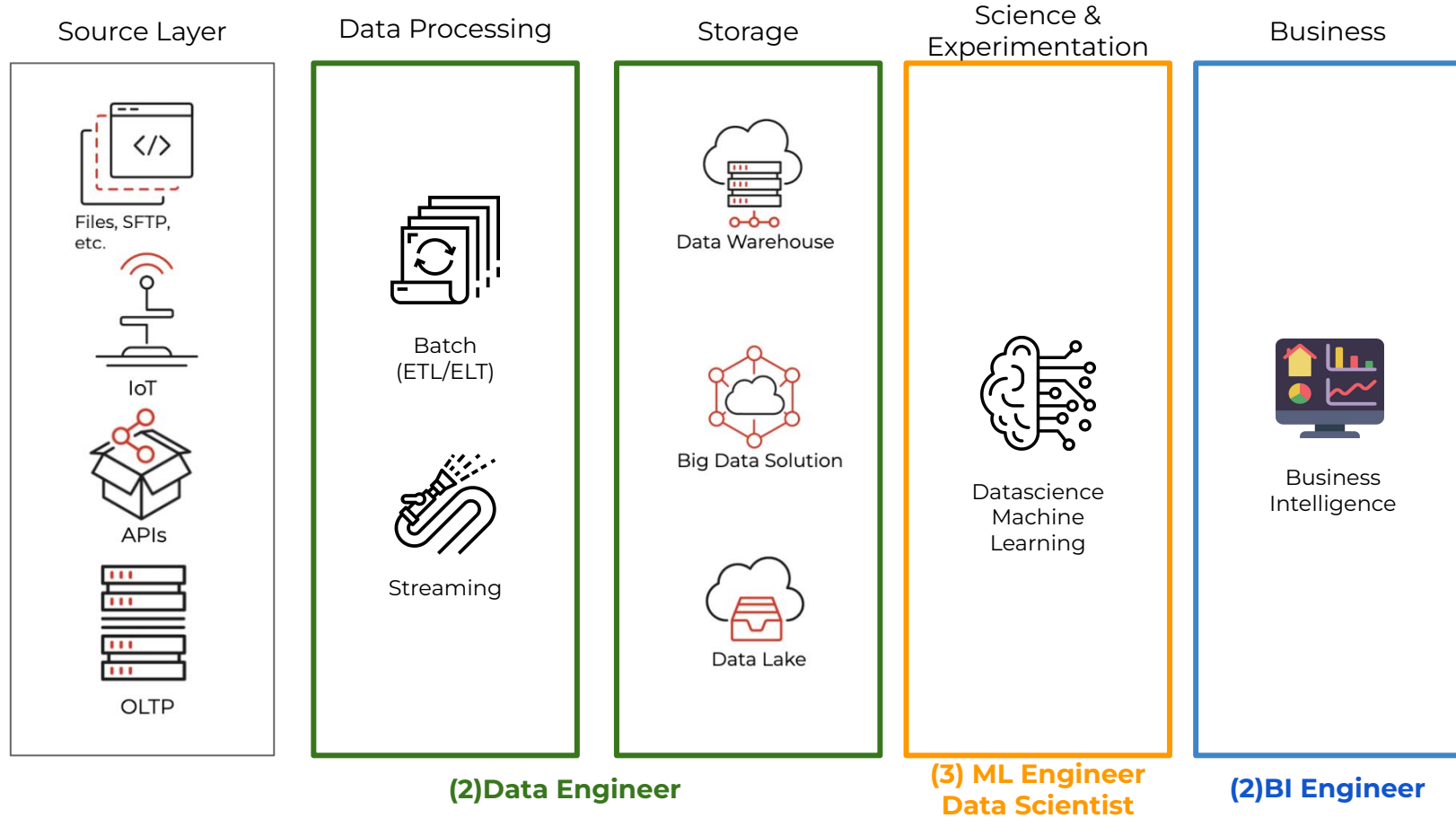


Business  
Intelligence



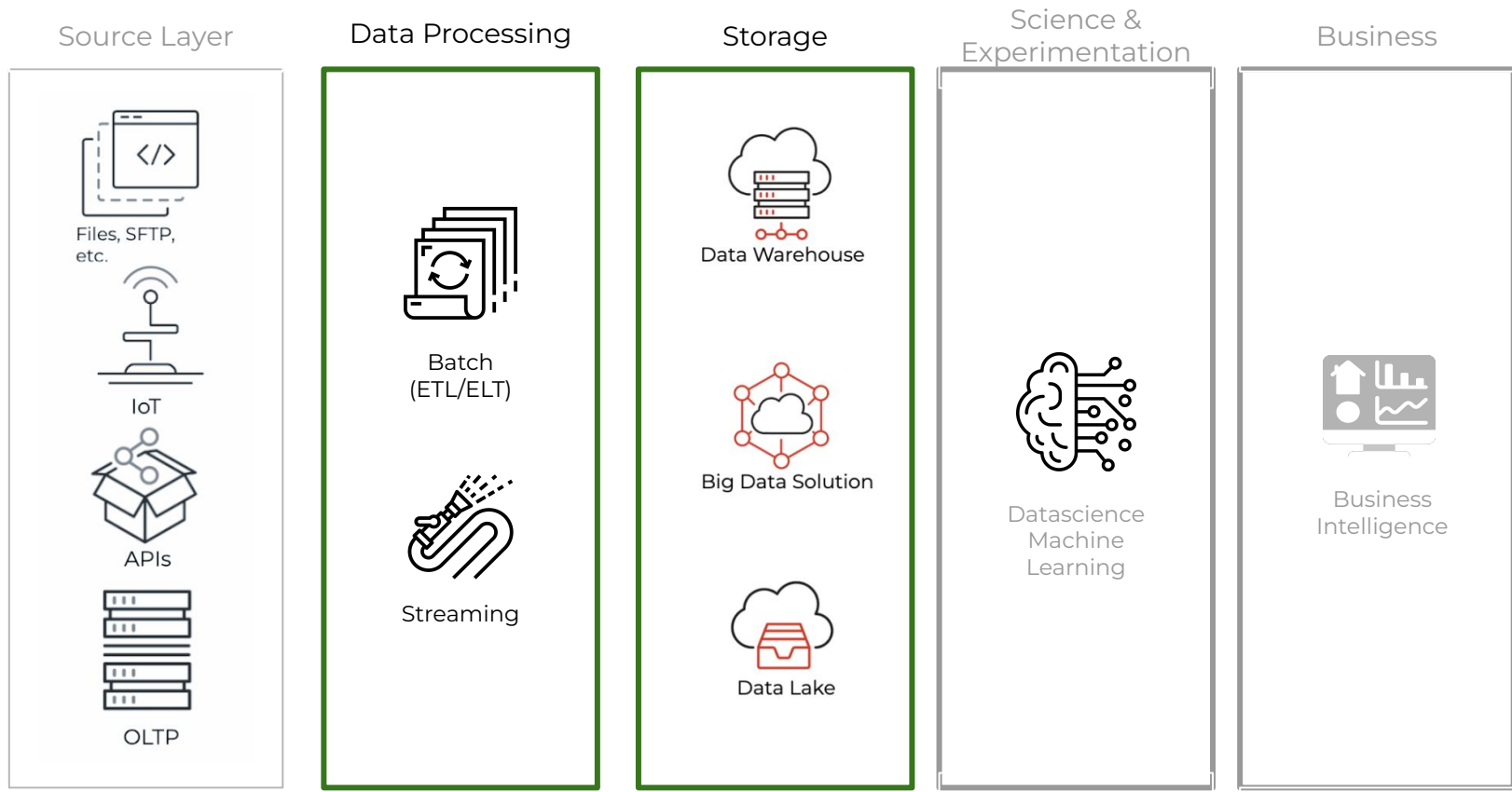
# Key Layers and roles

(1)Product Manager - manage data product.



# DE Key Layers

(1)Product Manager - manage data product.



(2)Data Engineer

(3) ML Engineer  
Data Scientist

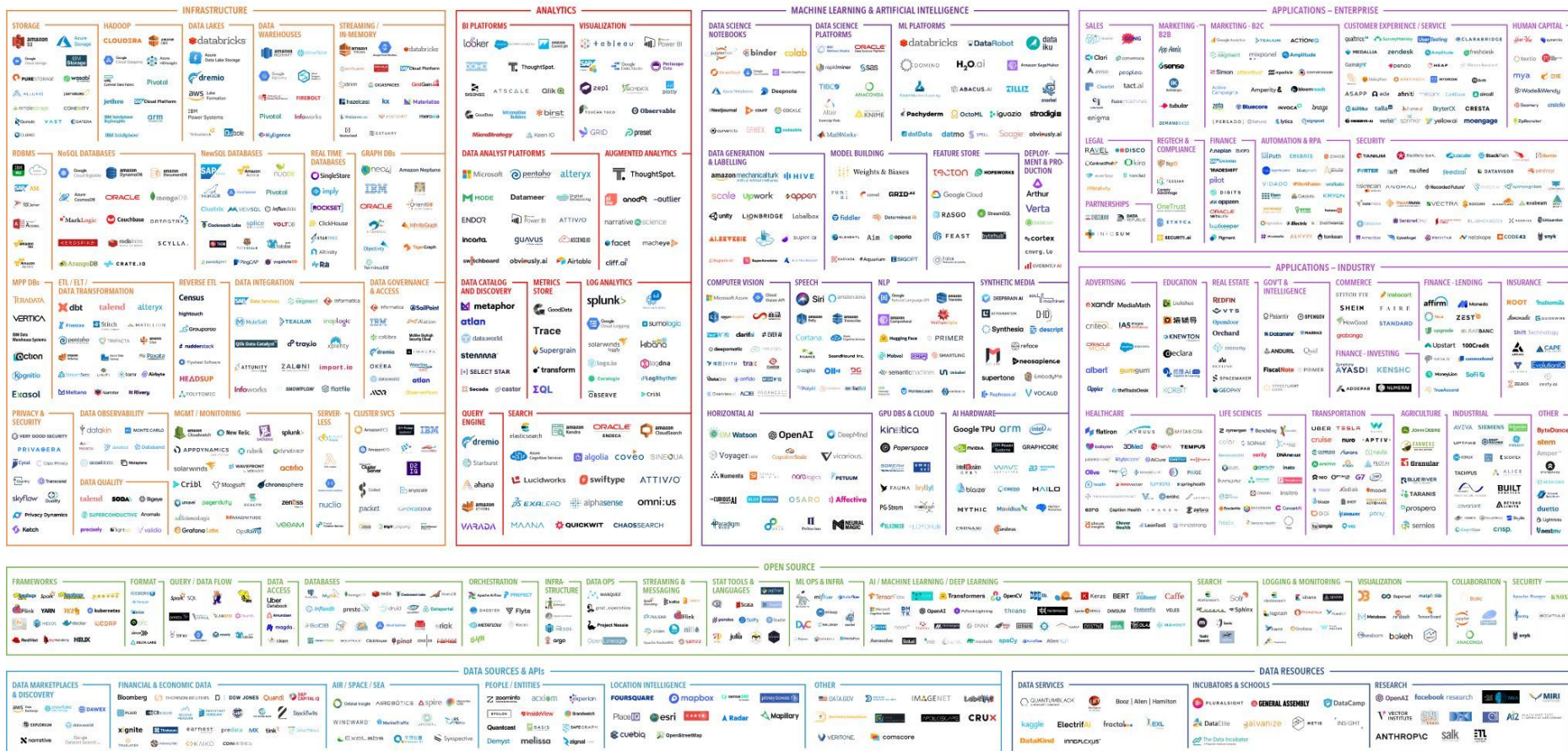
(2)BI Engineer

With DE you can - “move fast, break things”(c)...



# Data and AI Landscape 2021

MACHINE LEARNING, ARTIFICIAL INTELLIGENCE, AND DATA (MAD) LANDSCAPE 2021





# Data Stack with Open Source

## Source Layer



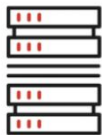
Files, SFTP,  
etc.



Game Client



APIs



OLTP

## Data Processing



## Storage



## Science & Experimentation



Spark Pool  
MLlib



## Business



# Data Stack for \$\$\$

## Source Layer



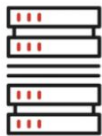
Files, SFTP,  
etc.



IoT



APIs



OLTP

## Data Processing

**Census**



**hightouch**

## Storage



**Google**  
Big Query



**amazon**  
REDSHIFT



**databricks**

## Science & Experimentation

**alteryx**



**DataRobot**

Datascience  
Machine  
Learning

## Business



**looker**

Business  
Intelligence

# Часть 2: Начало работы в Microsoft

# От отклика до начала работы

- Отклик Май
- 1 й звонок конец июля
- собеседования Август-Сентября
- Начало работы Ноябрь

Итоге: 5 месяцев





# Уровни Amazon vs Microsoft

Amazon	Microsoft
SDE I L4	SDE 59
	60
SDE II L5	SDE II 61
	62
SDE III L6 Senior SDE	Senior SDE 63
	64
Principal SDE L7	Principal SDE 65
	66
Senior Principal SDE L8	67
	Partner 68
Distinguished Engineer L10	69
	70
	Distinguished Engineer 80
	Technical Fellow

<https://www.levels.fyi/?compare=Amazon,Microsoft&track=Software%20Engineer>

# Что входит в оффер

- Sign Up bonus 1 и 2-й год
- RSU (акции на период 4 года)
- 10-20% годового бонуса
- базовая зарплата
- страховка
- отчисления на пенсию
- премиум переезд и 2 месяца в центре города для семьи

PS самый низкий total comp среди Faang, но хорошие бенефиты.

Filter + Add Compensation

Canada

Company Location   Date	Level Name Tag	Years of Experience Total / At Company	Total Compensation (CAD)		
			Base	Stock (yr)	Bonus
Microsoft Vancouver, BC, Canada   2 minutes ago	SDE II Data	12 yrs 2 yrs	\$185,500	159K   26.5K   N/A	
Microsoft London, ON, Canada   a day ago	Senior SDE Distributed Systems (Back-End)	10 yrs 3 yrs	\$192,495	150K   27.5K   15K	
Microsoft Vancouver, BC, Canada   a day ago	SDE II Web Development (Front-End)	5 yrs 0 yrs	\$191,141	121K   30.9K   39.2K	
Microsoft Toronto, ON, Canada   3 days ago	SDE II Data	5 yrs 5 yrs	\$251,750	198.8K   53K   N/A	
Microsoft Vancouver, BC, Canada   4 days ago	66 Distributed Systems (Back-End)	11 yrs 3 yrs	\$329,011	182.3K   100.3K   46.4K	
Microsoft Vancouver, BC, Canada   6 days ago	60 Web Development (Front-End)	1 yr 1 yr	\$141,200	107.2K   24K   10K	
Microsoft Vancouver, BC, Canada   7 days ago	SDE II Full Stack	15 yrs 2 yrs	\$144,613	110K   20.6K   14K	
Microsoft Vancouver, BC, Canada   10/07/2022	SDE API Development (Back-End)	2 yrs 1 yr	\$135,996	92K   24K   20K	
Microsoft Vancouver, BC, Canada   10/06/2022	66 Distributed Systems (Back-End)	10 yrs 10 yrs	\$456,338	190K   206.3K   60K	
Microsoft Vancouver, BC, Canada   10/06/2022	SDE II API Development (Back-End)	7 yrs 2 yrs	\$124,440	114.3K   9K   1.1K	

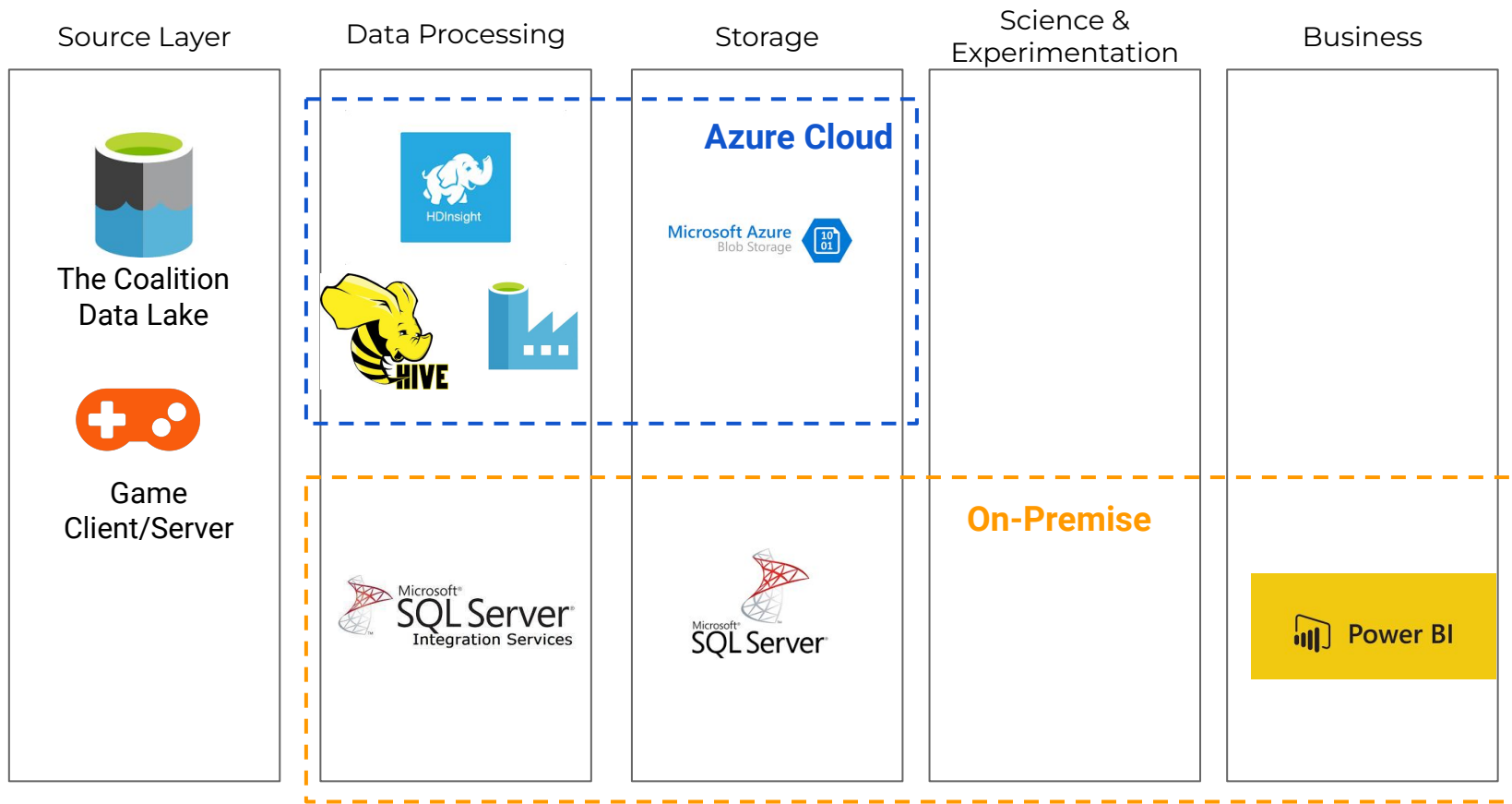
Rows Per Page  
10

1 - 10 of 248

1 2 3 4 ... 25 >

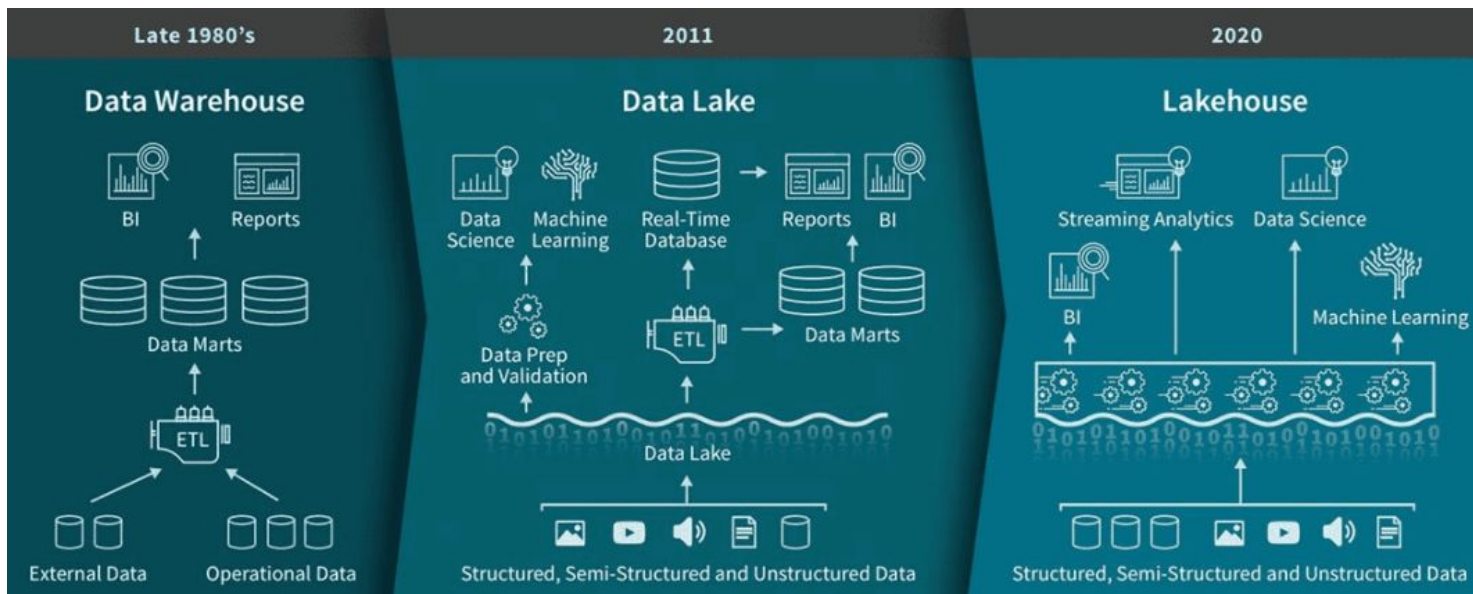
# Часть 3: Модернизация и миграция

# БЫЛО: HDInsights + Hive (staging), SQL Server SSIS (fact tables)



# Analytics architecture evolution

- До 2010 использовали хранилища данных (SMP, MPP).
- После появления Hadoop, все бросились строить озера данных. Преимущества - отдельный слой хранения и обработки данных, но нет поддержки ACID.
- Lake house = условно хранилище + озеро данных (Snowflake, Databricks, Synapse)

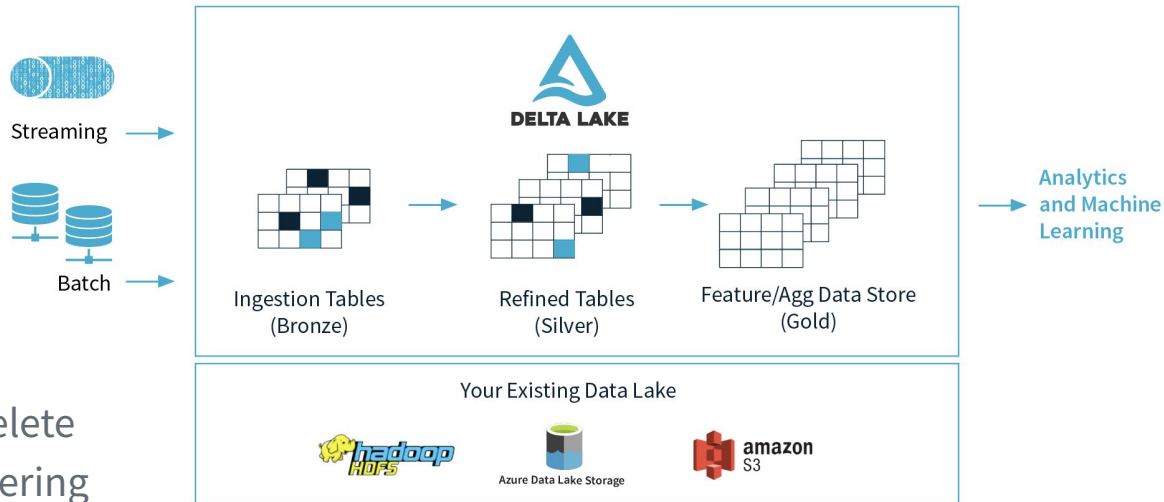


# Lake House = DW + Data Lake

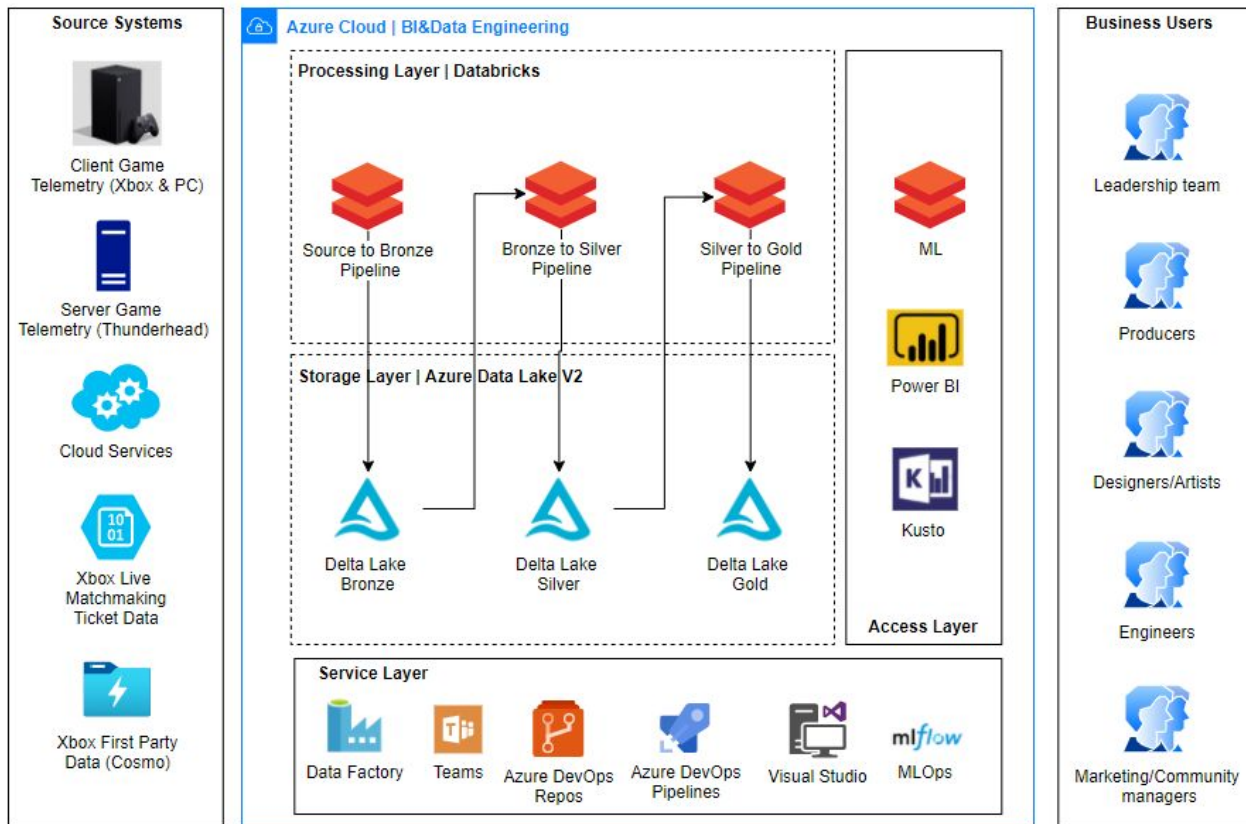
- Transaction Support (ACID)
- Schema Enforcement
- Upserts/Deletes

Key solutions for Lakehouse:

- **Apache Hudi** (Hadoop Update Delete and Incremental) by Uber Engineering
- **Apache Iceberg** by Netflix
- **Delta Lake** by Apache Spark



# Стало: Databricks+ Delta Lake



# Задача Инженера Данных в Gaming

- создание инфраструктуры данных для сбора телеметрии
- создание сред разработки, тестирования и продакшн
- трансформации бизнес логики в PySpark/SQL/C#
- автоматизация data pipelines
- качество данных, документация
- подключение к BI системам
- участие в разработки исходного формата и типа телеметрии
- тестирование игры

## Сложности:

- во время разработки (2-3) года у вас не ПРОД объема данных
- кол-во уникальных событий сотни, практически каждое событие имеет свою уникальную метрику



# Telemetry as a source of Player data

The word Telemetry is derived from the Greek roots tele, "remote", and metron, "measure".

Games are **state machines** - a person creates a **continual loop** of actions and responses which keep the game state changing. Often loops keeping the user **engages over a period of time**.

Telemetry helps to discovering **who** is performing **what** action **when** and **where** in the game. It cannot provide **why**.



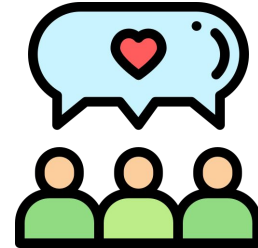
# 3 types of metrics



**Gameplay metrics**  
user behavior in the  
game



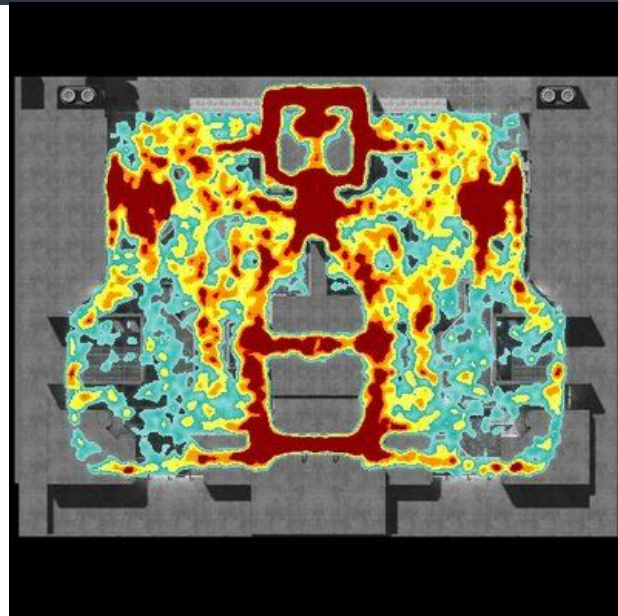
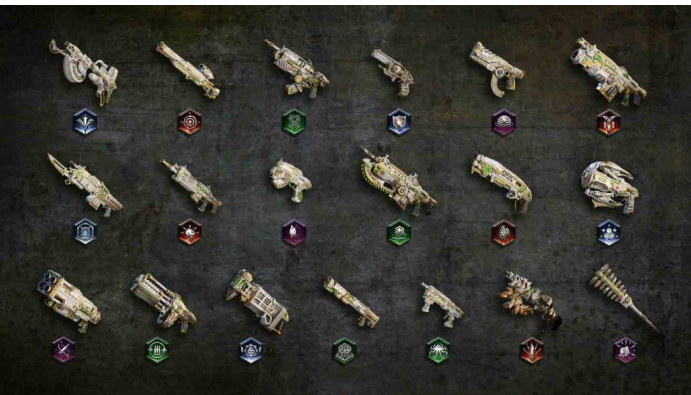
**Customer metrics**  
user as a customer,  
acquisition and  
retention



**Community metrics**  
user engagement in  
communities and social  
media

# Action Third-Person Shooters (TPS) Metrics

- Weapon use
- trajectory
- item/asset use
- character/kit choice
- level/map choice
- loss/win
- heatmaps
- team scores
- map lethality
- map balance
- vehicle use metrics
- special moves
- jumps and many more.



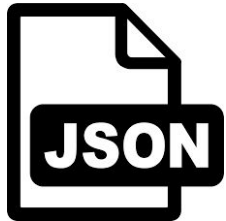
Death map, Halo3

# Data Engineering Design Flow as a Funnel



Event Names:

- Weapon Use
- Damage
- Shooting
- Flock
- Map Name
- HeartBeat
- and so on



Raw Tables (Bronze)  
Method: Append  
Trans: Minimum



Staging Tables (Silver)  
Method: Append  
Trans: JSON Schema



Fact Tables (Gold)  
Method: Merge  
Trans: Heavy

# Databricks custom logs

## Monitoring Azure Databricks in an Azure Log Analytics Workspace

---

This repository extends the core monitoring functionality of Azure Databricks to send streaming query event information to Azure Monitor. For more information about using this library to monitor Azure Databricks, see [Monitoring Azure Databricks](#)

The project has the following directory structure:

```
/src
  /spark-listeners-loganalytics
  /spark-listeners
  /pom.xml
/sample
  /spark-sample-job
/perftools
  /spark-sample-job
```

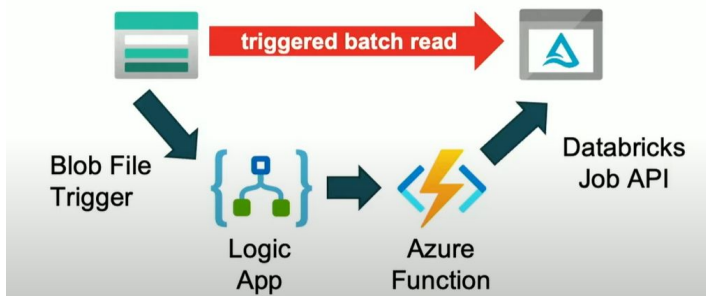
The **spark-listeners-loganalytics** and **spark-listeners** directories contain the code for building the two JAR files that are deployed to the Databricks cluster. The **spark-listeners** directory includes a **scripts** directory that contains a cluster node initialization script to copy the JAR files from a staging directory in the Azure Databricks file system to execution nodes. The **pom.xml** file is the main Maven project object model build file for the entire project.

The **spark-sample-job** directory is a sample Spark application demonstrating how to implement a Spark application metric counter.

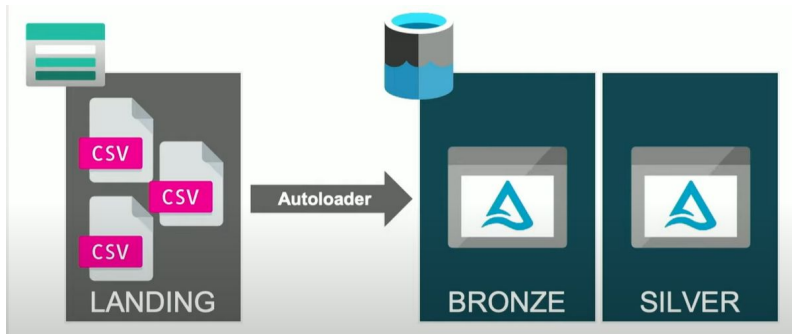
The **perftools** directory contains details on how to use Azure Monitor with Grafana to monitor Spark performance.

<https://github.com/mspnp/spark-monitoring>

# Databricks Delta Streaming (Auto Loader)



As files are discovered, their metadata is persisted in a scalable key-value store (RocksDB) in the checkpoint location of your Auto Loader pipeline. This key-value store ensures that data is processed exactly once.



## Delta Merge



```
def runThis(df, batchId):  
  (df  
   .write  
   .save(path)  
  )
```

```
df  
  .writeStream  
  .foreachBatch(runThis)  
  .save(path)
```



# Из интересного

Пожаловался лично Филу, что не дают работать инженеру и дают использовать сырой продукт Synapse...

Проблемы была решена оперативно.

## Phil Spencer named Microsoft Gaming CEO following Activision Blizzard deal

By Zachary Boddy published January 18, 2022

Microsoft's gaming leadership is uniting under Phil Spencer's guidance.



Comments (2)



Phil Spencer Xbox 20 Celebration (Image credit: Microsoft)

# Часть 4: Специфика DE работы в Microsoft



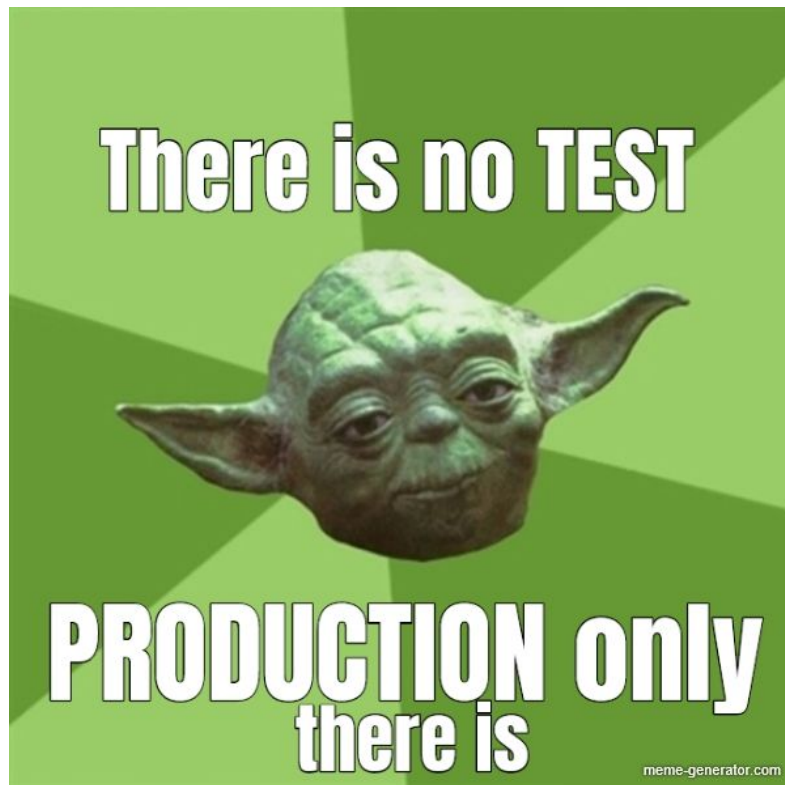
# Есть ли позиция Data Engineer в Microsoft?



# А кто есть?

- Software Development Engineer (SDE и DE)
- Data & Applied Scientist (analyst, BI, data science)
- Business Analyst (может быть и data analyst)

В Амазоне для ВІ и DE (SDE все Оки)



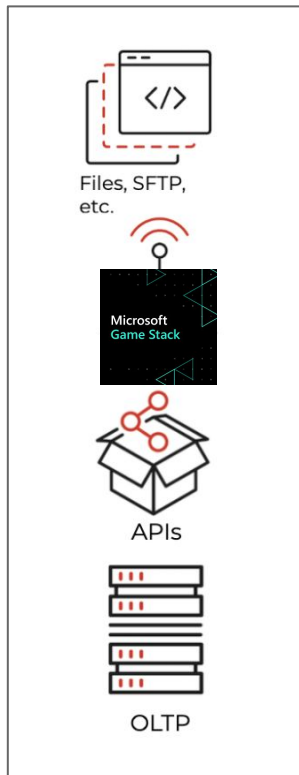
# В Microsoft нет DE роли, поэтому все Оки:)



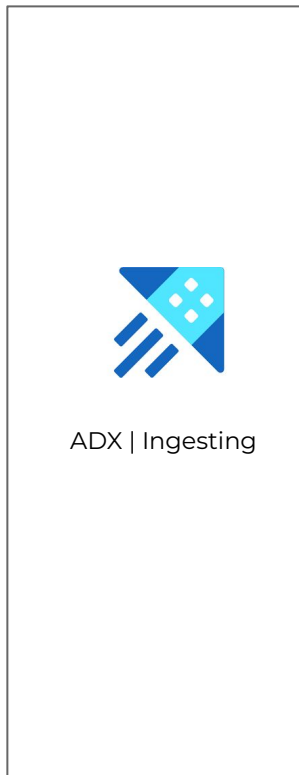
- Azure DevOps
- Visual Studio, Visual Code
- Yaml pipelines

# Типичное решение Azure Data Explorer (ADX)

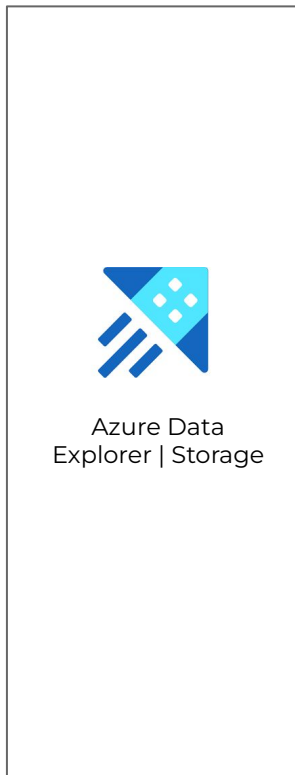
## Source Layer



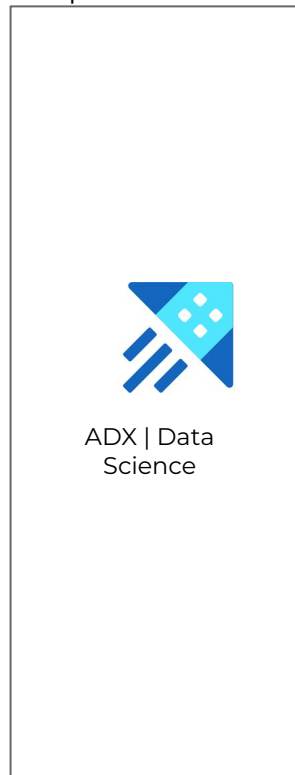
## Data Processing



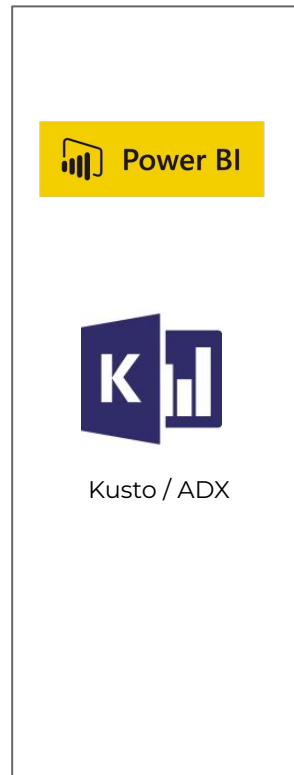
## Storage



## Science & Experimentation

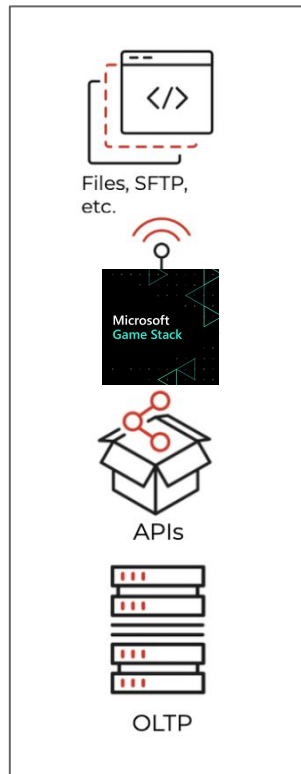


## Business

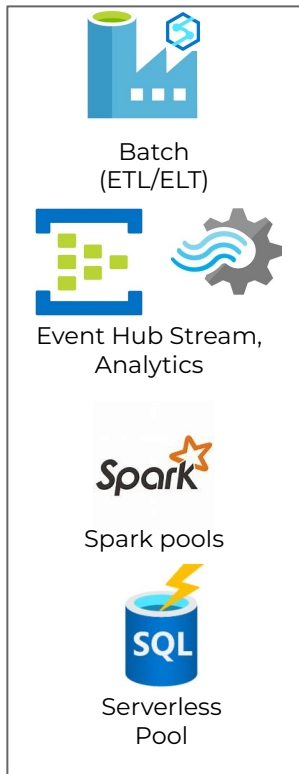


# Иногда такое Microsoft Azure Synapse

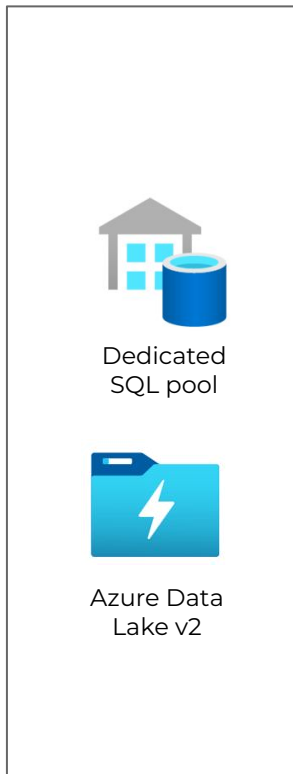
## Source Layer



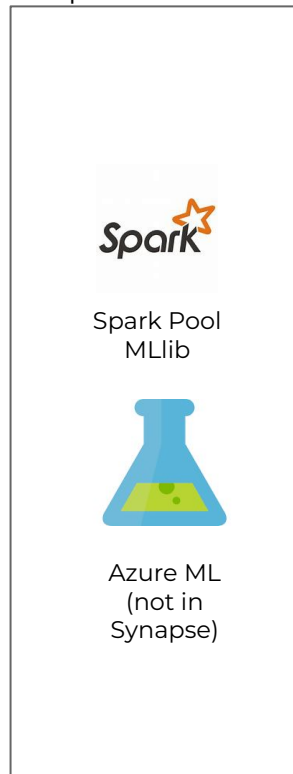
## Data Processing



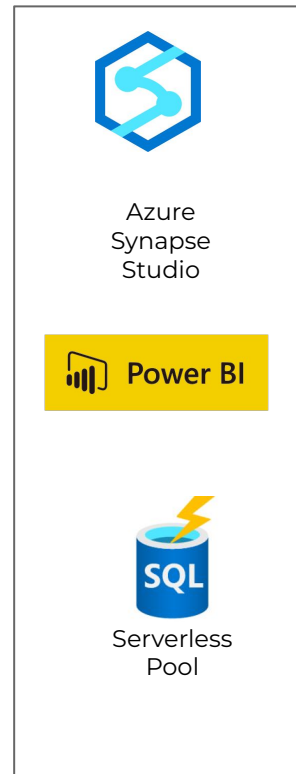
## Storage




## Science & Experimentation



## Business



# Microsoft Data Lake



Home > cosmos15-prod-cy > iebks.1pp.prod > shares > IEBKS.PartnerProd >

View ▾

Data Operations ▾

Options ▾

Add to favorites

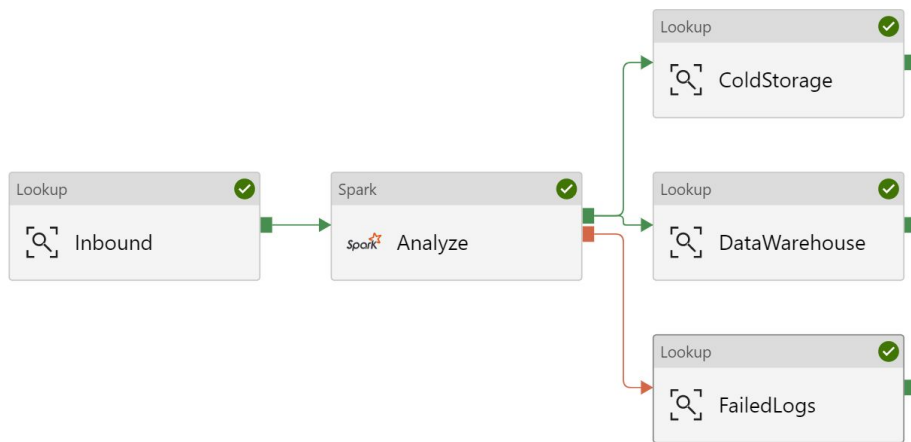
Settings ▾

<input type="checkbox"/>	Name	File size (Logical)	Storage consum
<input type="checkbox"/>		9.16 MiB	
<input type="checkbox"/>		7.37 MiB	
<input type="checkbox"/>		8.39 MiB	
<input type="checkbox"/>		7.36 MiB	
<input type="checkbox"/>		7.25 MiB	
<input type="checkbox"/>		7.26 MiB	
<input type="checkbox"/>		7.24 MiB	

## Scope script

```
test_input =  
  EXTRACT  
    FirstName : string,  
    LastName : string,  
    Age : int  
  FROM  
    "/local/Samples/SampleData/test_input.tsv"  
  USING DefaultTextExtractor();  
  
OUTPUT test_input  
  TO "/local/users/<your alias>/output.tsv"  
  USING DefaultTextOutputter();
```

# Data Factory Pipeline as a Code (C#)



```
ions.Generic;  
e.Management.DataFactory.Models;  
aFactory.Shared;
```

```
eDataFactory.Halifax
```

```
ass HalifaxPipeline : SharedPipeline
```

```
c HalifaxPipeline[] AllPipelines => new[]
```

```
Cook
```

```
new HalifaxPipeline
```

```
{  
    Name = "pl_Halifax_Cook",  
    PipelineResource = new PipelineResource
```

```
{  
    Folder = new PipelineFolder("Halifax"),  
    Parameters = new Dictionary<string, ParameterSpecification> { { "windowStart", new ParameterSpecification
```

```
Activities = new List<Activity>
```

```
{  
    //Pre-Processing
```

```
new GetMetadataActivity  
{
```

```
    Name = "GetMetadata-RtepClient-PreProcessing",  
    Dataset = new DatasetReference()
```

```
{  
    ReferenceName = "ds_Halifax_DataLakeStorage_Binary",  
    Parameters = new Dictionary<string, object> {  
        { "filePath", SharedSettings.DateHourFilePath("data-temp/rtep/client") },  
        { "fileName", new Expression("@coalesce(null)") }  
    }  
}
```

```
},
```



# Metadata ETL

Select Queries.sql...ulandrew.com (126) X

```
1 SELECT * FROM [procfwk].[ProcessingStageDetails]
2 SELECT * FROM [procfwk].[PipelineProcesses]
3 SELECT * FROM [procfwk].[PipelineParameters]
```

135 %

Results Messages

	StageId	StageName	StageDescription	Enabled
1	1	Extract	Ingest all data from source systems.	1
2	2	Transform	Transform ingested data and apply business logic.	1
3	3	Load	Load transformed data into semantic layer.	1

	PipelineId	StageId	PipelineName	Enabled
1	1	1	Stage 1-1	1
2	2	1	Stage 1-2	1
3	3	2	Stage 2-1	1
4	5	2	Stage 2-2	1
5	7	2	Stage 2-3	1
6	9	3	Stage 3-1	1

	ParameterId	PipelineId	ParameterName	ParameterValue
1	1	9	TestParam1	Frank
2	2	9	TestParam2	Bob

LocalExecutionId
StageId
PipelineId
PipelineName
StartDateTime
PipelineStatus
EndDateTime

LogId
LocalExecutionId
StageId
PipelineId
PipelineName
StartDateTime
PipelineStatus
EndDateTime

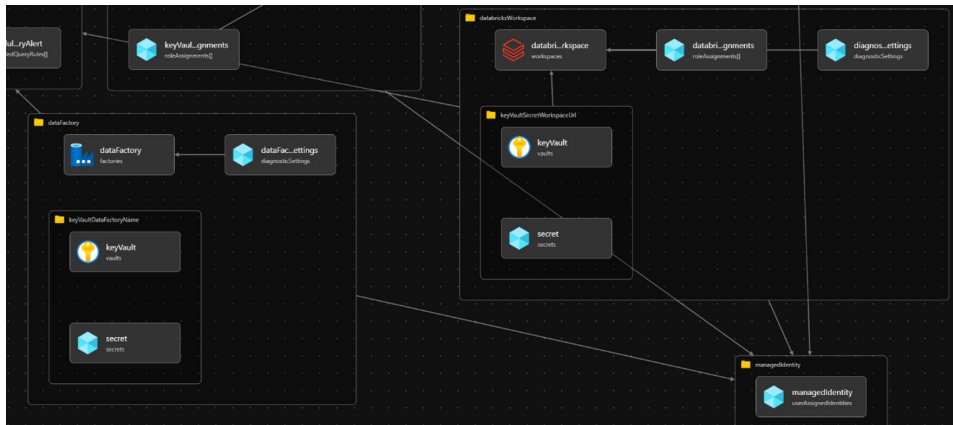
PipelineId
StageId
PipelineName
Enabled

StageId
StageName
StageDescription
Enabled

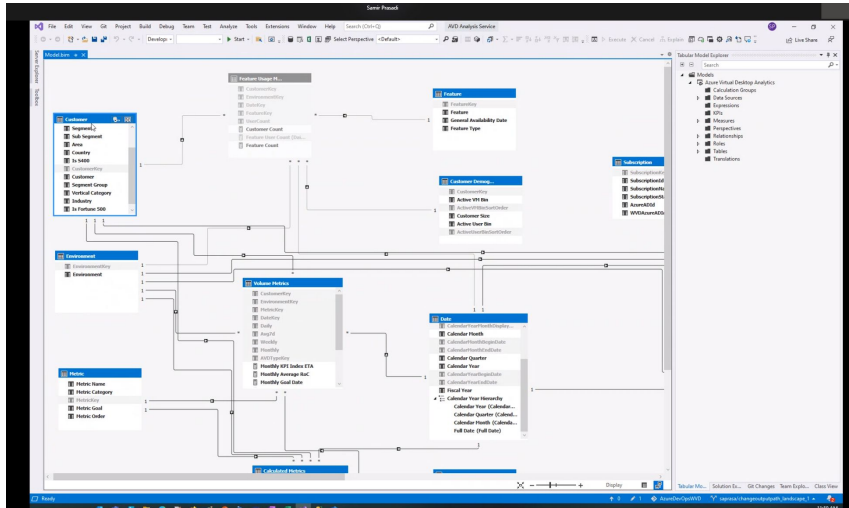
ParameterId
PipelineId
ParameterName
ParameterValue

# Infrastructure as a Code – Azure Bicep

```
databricks.bicep X
modules > databricks.bicep > ...
1  targetScope = 'resourceGroup'
2
3  @description('Location for all resources.')
4  param location string
5
6  @description('The name of the network security group to create.')
7  param workspaceName string
8
9  var managedResourceGroupName = 'databricks-rg-${workspaceName}-${uniqueString(workspaceName, resourceGroup().id)}'
10
11 resource databricksWorkspace 'Microsoft.Databricks/workspaces@2021-04-01-preview' = {
12   name: workspaceName
13   location: location
14   properties: {
15     managedResourceGroupId: subscriptionResourceId('Microsoft.Resources/resourceGroups', managedResourceGroupName)
16   }
17 }
18
19 output workspaceUrl string = databricksWorkspace.properties.workspaceUrl
20
```



# BI as a Code

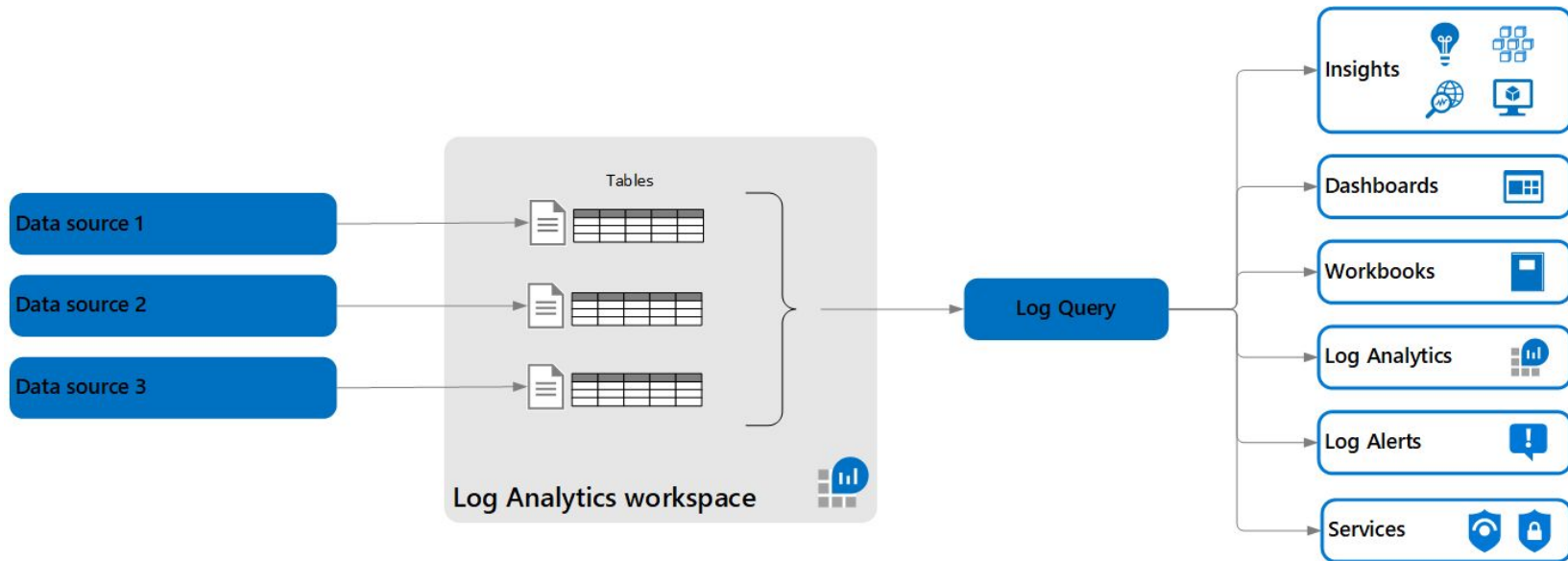


The screenshot displays the AVD Studio interface with a JSON configuration for a data source. The configuration is for a 'Server' named 'sigmaadbs.database.windows.net' and a 'Database' named 'RDSD'. The configuration includes authentication details, a query, and a list of tables and columns. The 'Tables' section lists 'Customer' and 'Sales' tables, and the 'Columns' section lists various attributes like 'Region', 'Country', and 'Product'. The configuration is as follows:

```
13 {
14   "server": "sigmaadbs.database.windows.net",
15   "database": "RDSD",
16   "authentication": null,
17   "query": null,
18   "credential": {
19     "authenticationKind": "UsernameAndPassword",
20     "kind": "SQL",
21     "path": "sigmaadbs.database.windows.net/RDSD",
22     "username": "sigmaadbsanalysis@sigmaadbs",
23     "password": "sigmaadbsanalysis@sigmaadbs",
24     "encryptConnection": true
25   },
26   "tables": [
27     {
28       "name": "Customer",
29       "columns": [
30         {
31           "name": "AzureADId",
32           "dataType": "string",
33           "sourceColumn": "AzureADId"
34         },
35         {
36           "name": "IPId",
37           "dataType": "string",
38           "description": "Data team based ID that is assigned to an organization.",
39           "isNullable": false,
40           "sourceColumn": "IPId"
41         },
42         {
43           "name": "Vertical",
44           "dataType": "string",
45           "sourceColumn": "Vertical"
46         },
47         {
48           "name": "Segment",
49           "dataType": "string",
50           "sourceColumn": "Segment"
51         },
52         {
53           "name": "Sub Segment",
54           "dataType": "string",
55           "sourceColumn": "SubSegment"
56         },
57         {
58           "name": "Area",
59           "dataType": "string",
60           "sourceColumn": "Area"
61         }
62       ]
63     }
64   ]
65 }
```

# Alerting

Azure Log Analytics (по сути ADX) умеет собирать данных из всех сервисов Azure и позволяет создавать нам Alerts на ошибки с использование Query. Также можно создавать свои собственные Alerts.



# Мои выводы

- У Microsoft своя экосистема
- Нужно любить продукты Microsoft
- Отличная интеграция между продуктами
- Часто приходится работать с сырыми продуктами (DOGFOOD)
- Не самые высокие зарплаты, но хорошие бенефиты
- В Microsoft хорошо попадать “молодым” специалистам в начале карьеры.
- Отличный опыт как делать решения на Azure