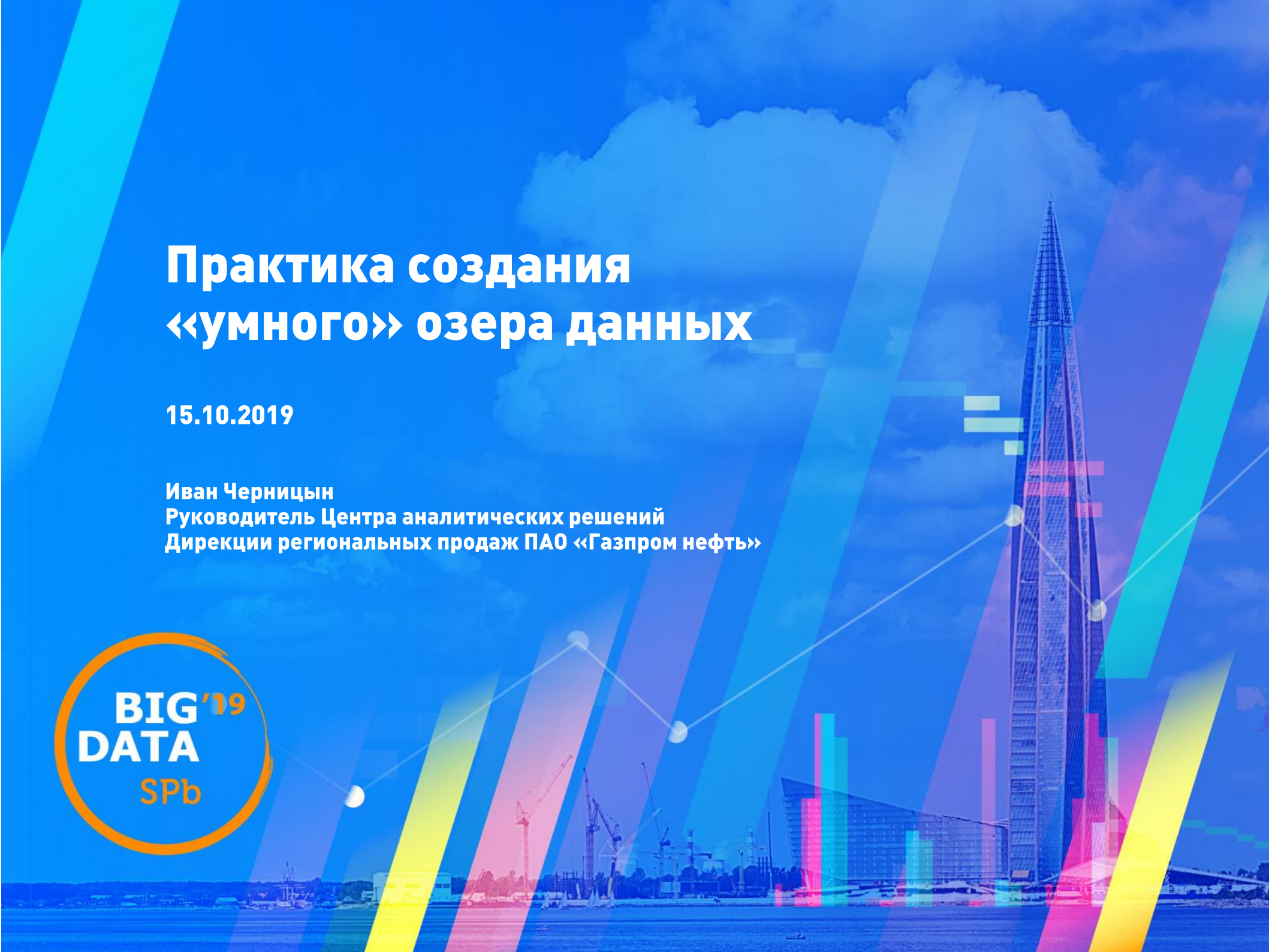


Практика создания «умного» озера данных

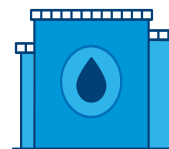
15.10.2019

Иван Черницын
Руководитель Центра аналитических решений
Дирекции региональных продаж ПАО «Газпром нефть»

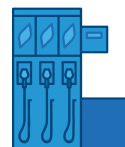
BIG^{'19}
DATA
SPb



- Биржевые продажи нефтепродуктов
- Реализация нефтепродуктов крупным и мелким оптом
- Хранение нефтепродуктов
- Управление >50 собств. нефтебазами



- Розничная реализация нефтепродуктов
- Продажи корпоративным клиентам
- Управление АЗС \ АЗК >1800 объектов



- >800 розничных магазинов и кафе при АЗС
- Оказание услуг моек, СТО



- Управление автотранспортом: бензовозы/газовозы
- Услуги метрологии



29 регионов присутствия в России + 4 страны СНГ

> 20 тысяч сотрудников

Клиентов-участников бонусной программы > 11,4 млн. человек

Стандартизация бизнес-процессов, оргструктур, внедрение унифицированных шаблонных систем породили первые практики управления данными и потребность в аналитических решениях

- Система управления нормативно-справочной информацией
- Система управления АЗС
- Система коммерческого учета
- Система бухгалтерского и управленческого учета
- Система планирования и бюджетирования
- Система управления инвестициями
- Система управления персоналом
- Система управления МТО
- Система управления ТО и ремонтами
- Система управления автотранспортом



**Централизация
методологии**

**Институт владельцев
систем**

**Единые ключевые
справочники**

**Высокая потребность в
интеграции данных**



*Высокий темп организационных изменений –
существенная реорганизация каждые 2 года*





> 1 100
пользователей



>90 систем-источников
>50 внешних и несистемных
источников данных



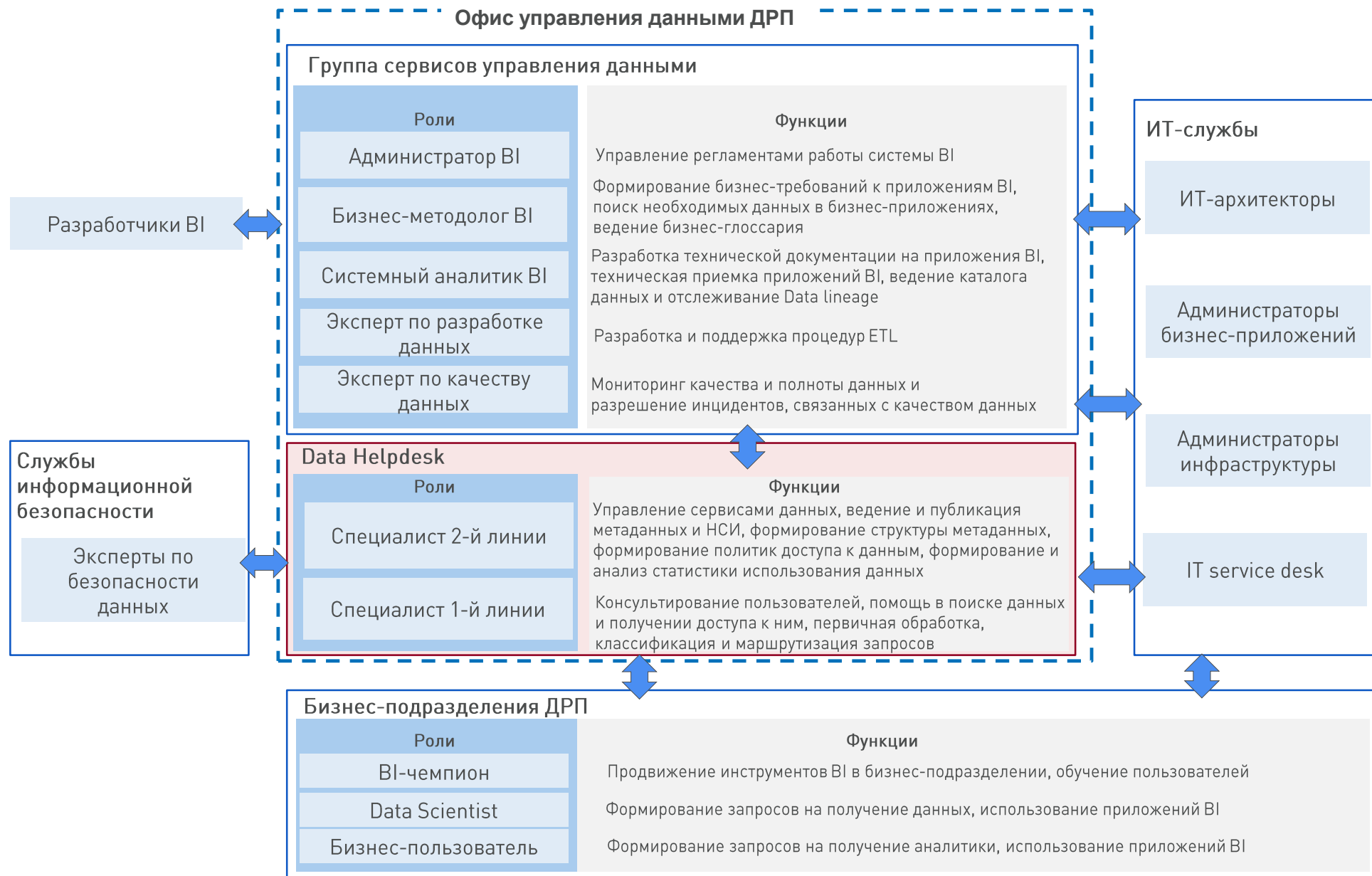
> 250
аналитических
приложений

Функциональные подсистемы BI охватили ключевые направления деятельности:

- Учет и консолидация
- Экономический анализ
- Расходы
- Финансы
- Инвестиции
- Маржинальный доход
- Материальный баланс
- АРМы руководителей

- Аналитика розничных продаж
- Клиентская аналитика
- Сервисы для АЗС
- Сравнительный анализ АЗС
- Программа лояльности АЗС
- Управление ассортиментом
- Анализ работы оборудования
- Мониторинг цен

- ИТ-сервисы
- Управление проектами
- Аналитика персонала
- Эффективность бизнес-процессов
- Оптовые продажи
- Логистика
- Качество данных



Making Your Landscape Less Uncertain — Focus Areas



Component —
not core



Operational and
analytics skills



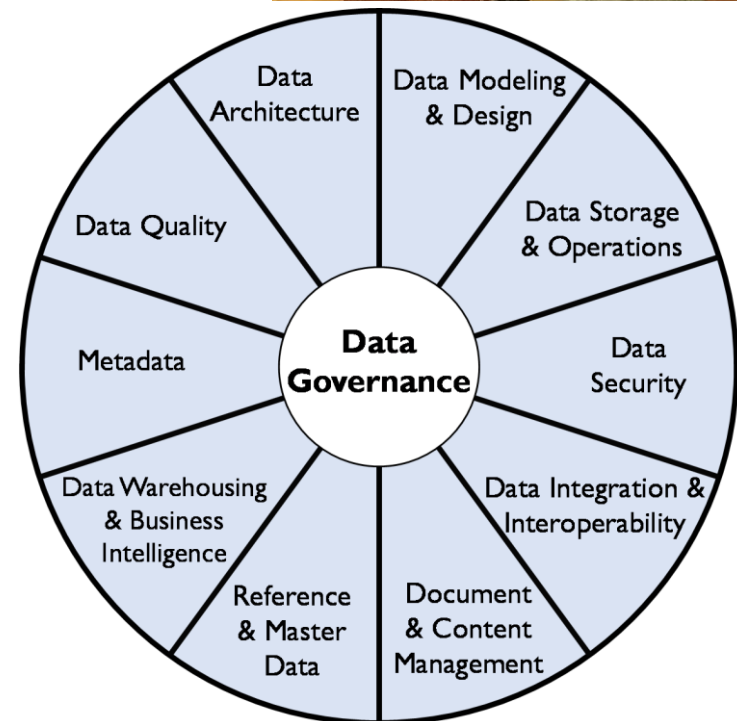
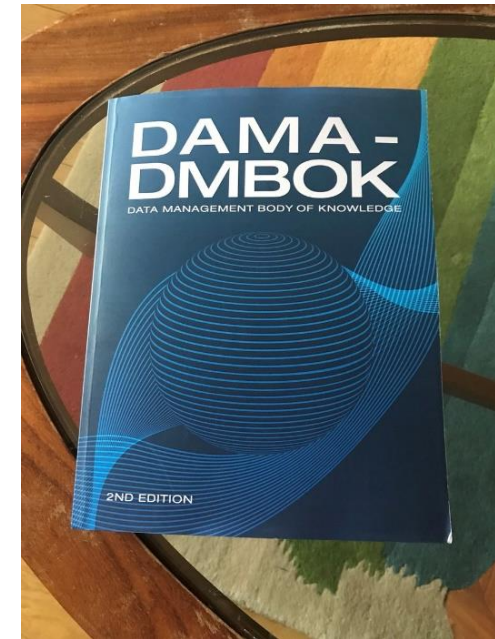
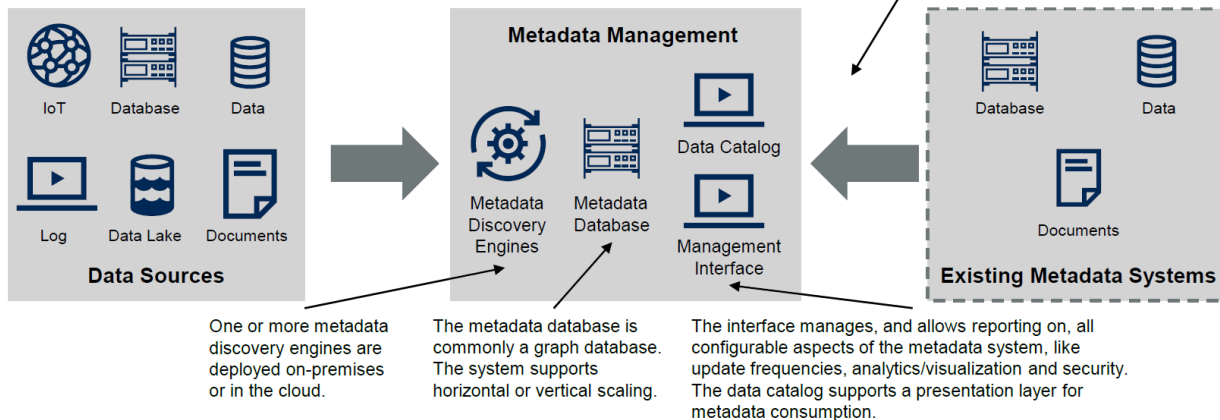
Security and
metadata

Data Catalogs Are the New Black in Data Management and Analytics

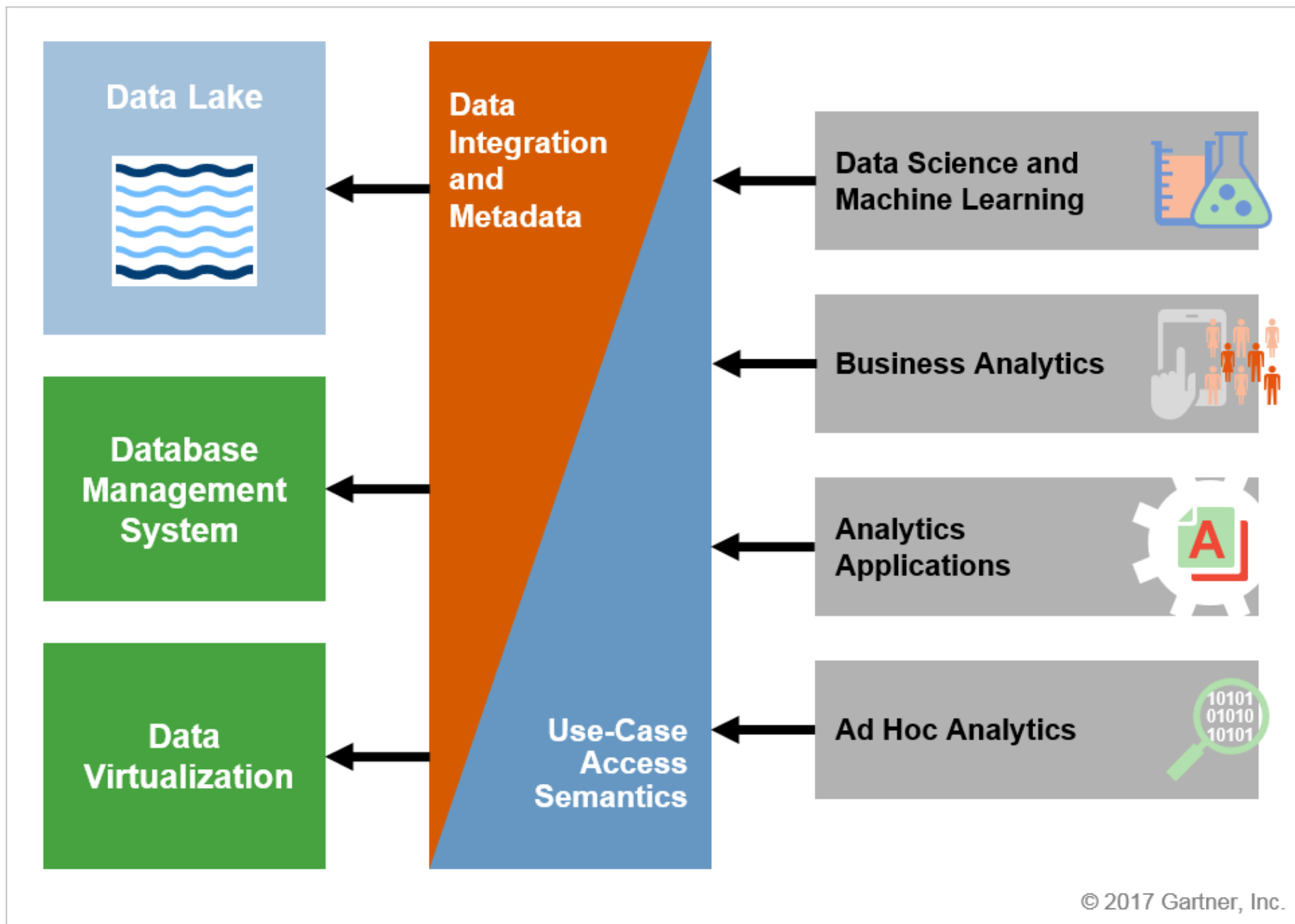
Published: 13 December 2017 ID: G00338777

Analyst(s): Ehtisham Zaidi, Guido De Simoni, Roxane Edjlali, Alan D. Duncan

Metadata Management Landscape



1. **«Озеро данных»** - промышленные инструменты интеграции, ETL и хранения данных
2. **Единый каталог правил по качеству данных** с управлением и мониторингом
3. **Система управления метаданными** для аналитиков: каталог доступных источников, таблиц, витрин, моделей данных и отчетов + цепочки происхождения
4. **Пользовательский портал по данным**: поиск определений данных, владельцев, ролей, отчетов и т.д.






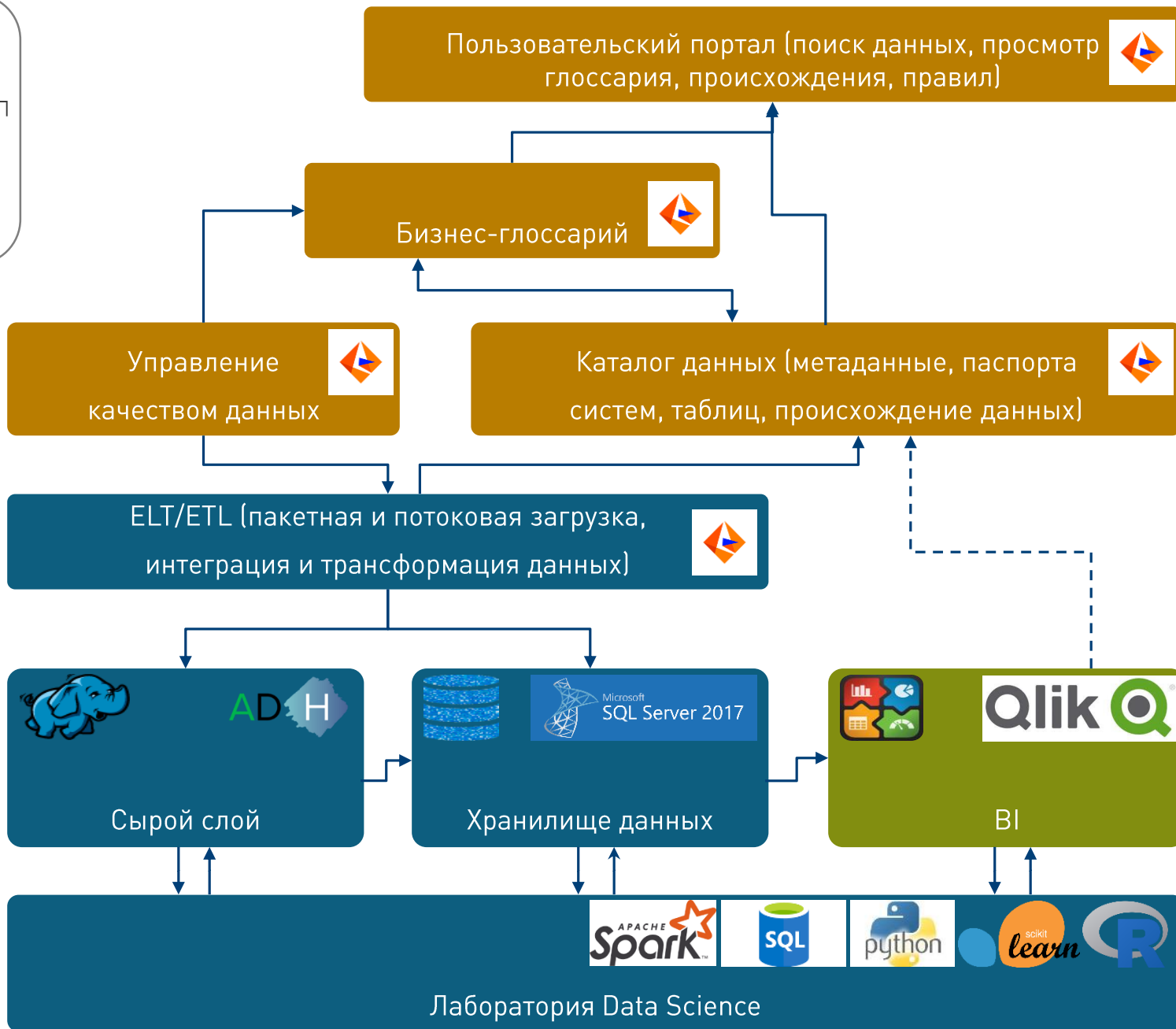
© 2017 Gartner, Inc.



- Ведение бизнес-метаданных и инструменты управления изменениями
- Бизнес-глоссарий с описанием бизнес-терминов, их взаимосвязей и атрибутов, связанных правил по качеству данных, процессами управления изменениями
- Бизнес-глоссарий, **интегрированный с подсистемой управления качеством данных**
- В контуре «озера данных» **сбор потоков происхождения данных** от систем источников до конечных пользовательских представлений и BI-приложений **производится в автоматизированном режиме**
- В подсистеме управления метаданными реализован мониторинг изменений технических метаданных
- Единый каталог правил по качеству данных
- **Подсистема управления качеством данных, интегрированная с «озером данных»**
- Измерение, отображение и аудит всех проверок по качеству данных в контуре «озера данных» **в разрезе потоков обработки данных**
- Паспортизация источников и полей данных
- Функциональность профилирования данных
- Функциональность разметки данных (коммерческая тайна и персональные данные)
- Подсистема для разрешения запросов бизнес-пользователей на поиск доступных данных и ролей для доступа к данным, **интегрированная с бизнес-глоссарием и системой управления метаданными** | 12

Легенда

-  Компоненты Системы управления данными ДРП
-  Компоненты «Озера данных» ДРП
-  Компоненты платформы Informatica





Сегменты Клиентского контура

Базовая сегментация:
АВТО, КОМ ТРАНС, В2G, НАЗС, ПРОМ СХ

Подсегменты: Авто, Fleet, CRT, НКП, ПРОМ, СХ

Активность: активные, неактивные, архив

Новизна клиента

Лояльность G-Drive

Вид деятельности по ОКВЭД

Объем потребления клиента:
массовый, средний, крупный

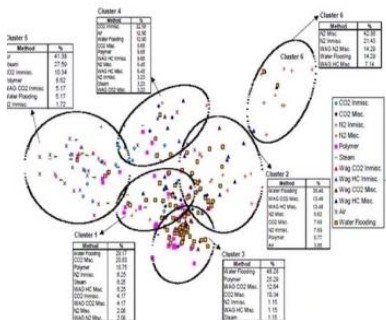
Активность СТиУ

Тип / класс / износ ТС

Возраст

Пол

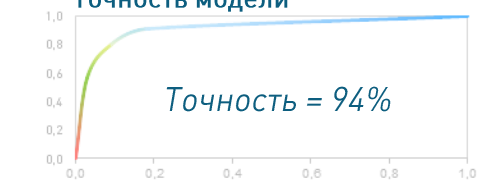
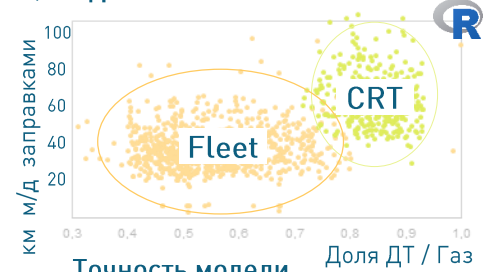
География



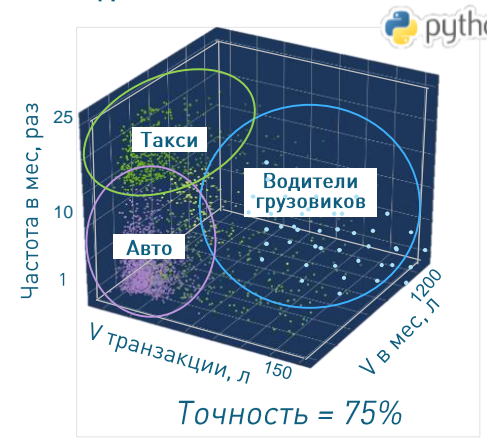
- Используются методы для сегментаций:
- Classification Tree
 - Random Forest
 - XGB
 - GBM
 - Нейронные сети

Сегментация клиентов на основе методов машинного обучения и продвинутой аналитики

1) Подсегменты КП



2) Подсегменты ПЛ



- ✓ С мая 2019 года все новые кейсы ДРП реализуются по целевому сценарию: разработка ETL на продуктах Informatica с автоматической передачей трансформаций в Систему управления данными

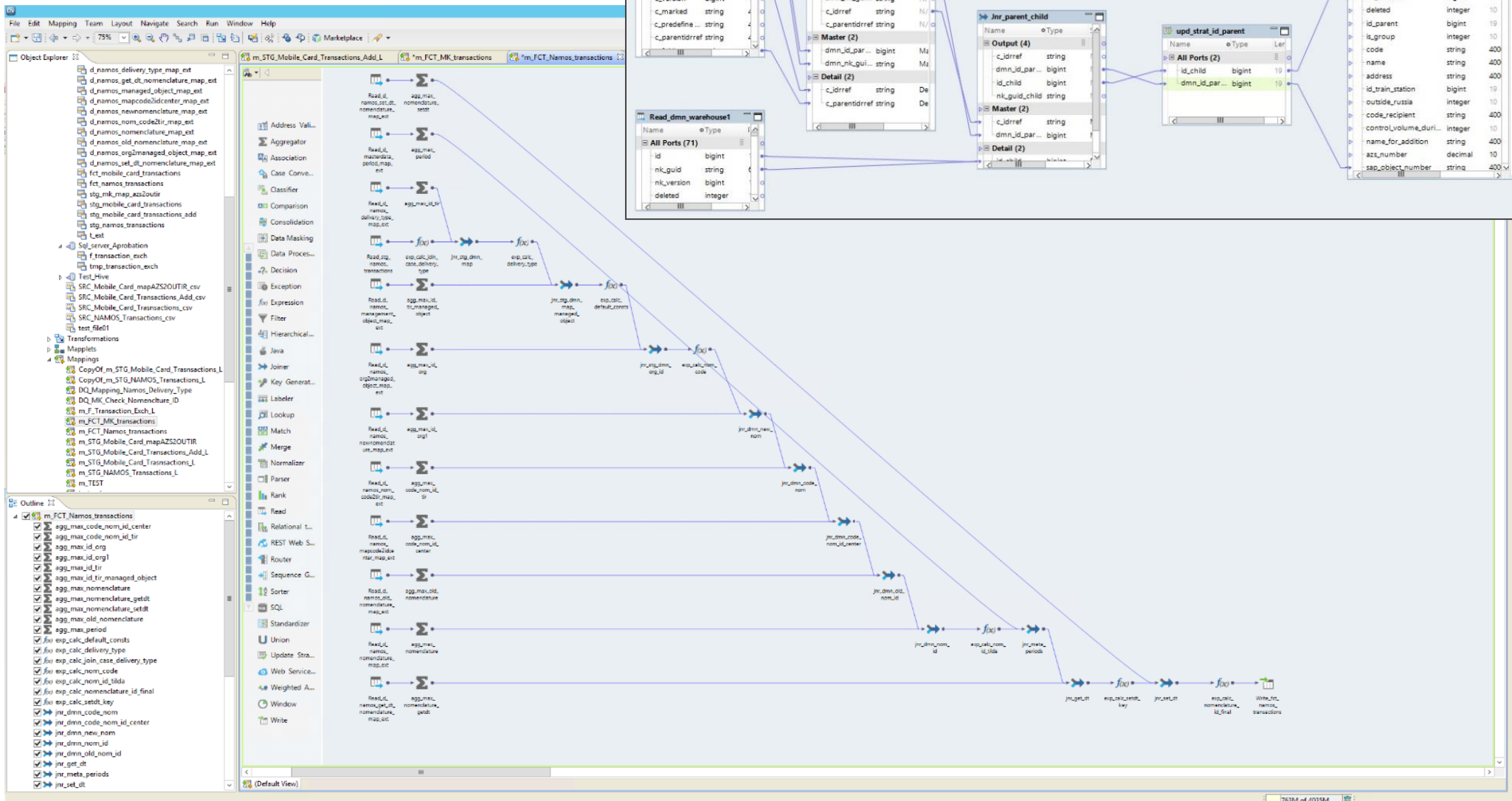
Далее на слайдах примеры реализованного функционала



- Особенности:
- Требуется изменение культуры разработки
 - Применяем итерационный подход к поиску лучших практик
 - Вырабатываем шаблоны для ускорения разработки



! Требуется значительное количество соглашений по практике разработки для автоматической работы инструментов управления данными



Подсистема качества данных интегрирована с инструментами ETL и позволяет создавать и встраивать правила по качеству данных в ETL-процессы

Informatica | New | Open | Search All

Start | Glossary | **Discovery** | Design | Scorecards

DQ_NAMOS_LAYER0

1 Specify General Properties | 2 Select Source | 3 Specify Settings | **4 Specify Rules**

7 Applied

Name
DQ_Namos_PLOTNOST_00001
DQ_Namos_VREMYAOPERACII_00001
DQ_Namos_VREMYAOPERACII_00002
DQ_NAMOS_RNN_00001

DQ_Namos_VREMYAOPERACII_00001

Name: DQ_Namos_VREMYAOPERACII_00001

Rule Condition:: IIF ((DATE_COMPARE(DATA_I_VREMYA_NACHALA_SMENY, DATA_I_VREMYA_REGISTRACII_PRODAZHI)=1) , '0', '1')

Output Column: DQ_Namos_VREMYAOPERACII_00001

DQ_Namos_VREMYAOPERACII_00002

Name: DQ_Namos_VREMYAOPERACII_00002

Rule Condition:: IIF ((DATE_COMPARE(DATA_I_VREMYA_REGISTRACII_PRODAZHI, DATA_I_VREMYA_OKONCHANIYA_SMENY)=1) , '0', '1')

Output Column: DQ_Namos_VREMYAOPERACII_00002

DQ_NAMOS_RNN_00001

Name: DQ_NAMOS_RNN_00001

Rule Condition:: IIF(RRN != RRN_MARKETINGOVOJ_TRANZAKCII, '1', '0')

Output Column: DQ_NAMOS_RNN_00001

Informatica | New | Open | Search All | Manage | ?

Start | Glossary | **Discovery** | Design | Scorecards

DQ_NAMOS_LAYER0

Back to Overview | Profile Run 5 of 5 | 59 of 62 Columns | 7 of 7 Rules | 100000 Rows | 08/14/2019 11:34:06 AM

Columns and Rules (20) | DQ_NAMOS_RNN_00001 | < > | DQ_NAMOS_RNN_00001

General

Total rows 100000

Null: 0 (0%)

Distinct: 2 (0.01%)

Non-distinct: 99998 (99.99%)

fx Rule used: DQ_NAMOS_RNN_00001
Input columns: RRN, RRN_MARKETINGOVOJ_TRANZAKCII

Add Tag | Add Comment

Datatype

Datatype	Percentage	Documented
string(10)		
Decimal(1)	100.00%	
Fixed Length String(1)	100.00%	
Integer(1)	100.00%	
String(1)	100.00%	

Patterns

Pattern	Percentage	Documented
9	100.00%	

Values

2 distinct values (2 Non-unique, 0 Unique) | Sort By: Frequency | Descending

Value	Frequency	Length	Percentage
1	81389	1	81.39%
0	18611	1	18.61%

Length (min → max) 1 → 1
Value (min → max) 0 → 1
Average: 0 | Sum: | Standard Deviation: 0

Data Preview

[DQ_NAMOS_RNN_00001 : "0"]
First 100 rows only

	IDENTIFIKATO...	NOMER_SMENY	DATA_I_VREMY...	DATA_I_VREMY...	DATA
1	'0000025061	920	10/07/2019 04:0...	10/08/2019 03:5...	10/07/...
2	'0000025061	920	10/07/2019 04:0...	10/08/2019 03:5...	10/07/...
3	'0000025061	920	10/07/2019 04:0...	10/08/2019 03:5...	10/07/...

Доступно управление типами (аспектами) и уровнем критичности правил

DS-159: fct_stock_index
Data Set in DW_SPB_SPB99_DBT01_DWH_DEV03

SUMMARY ATTRIBUTES RELATIONSHIPS STAKEHOLDERS VALUES DATA QUALITY IMPACT HISTORY CHANGE FOLLOW

Dashboard Rules

DETAILS

Cont...	Rule Type	Ref.	Description	A...	Attributes	System	Data Set	Target	Last ...	Last Result Date
8 records										
	Допустимость	⊙ DQ_INTERNET_StockIndexInde	Числовое значение не отрицательное.	Active	value_wo_trar	DW_SPB_SPB99_DB'	fct_stock_index	90%	100%	24-Jun-2019
	Допустимость	⊙ DQ_INTERNET_StockIndexNum	Числовое значение не отрицательное	Active	num_trades	DW_SPB_SPB99_DB'	fct_stock_index	90%	100%	12-Aug-2019
	Допустимость	⊙ DQ_INTERNET_StockIndexDate	Формат не соответствует маске DD.MM.YYYY	Active	dt	DW_SPB_SPB99_DB'	fct_stock_index	90%	100%	13-Aug-2019
	Допустимость	⊙ DQ_INTERNET_StockIndexInde	Числовое значение не отрицательное	Active	index_value	DW_SPB_SPB99_DB'	fct_stock_index	90%	100%	26-Jul-2019
	Допустимость	⊙ DQ_INTERNET_StockIndexValu	Числовое значение не отрицательное	Active	value_per_lite	DW_SPB_SPB99_DB'	fct_stock_index	90%	100%	30-Jul-2019
	Допустимость	⊙ DQ_INTERNET_StockIndexValu	Числовое значение не отрицательное	Active	value_per_lite	DW_SPB_SPB99_DB'	fct_stock_index	90%	100%	26-Jul-2019
	Допустимость	⊙ DQ_INTERNET_StockIndexValu	Числовое значение не отрицательное	Active	value	DW_SPB_SPB99_DB'	fct_stock_index	90%	100%	26-Jul-2019
	Допустимость	⊙ DQ_INTERNET_StockIndexVolu	Числовое значение не отрицательное	Active	volume	DW_SPB_SPB99_DB'	fct_stock_index	90%	100%	26-Jul-2019

8 records

Informatica New Open

Start Glossary Discovery Design Scorecards Library

DQ_NAMOS_LAYER0

Columns and Rules (20) DQ_NAMOS_RIN_00001

General

Total rows 100000

2 distinct values (2 Non-unique, 0 Unique)

Null 0 0%

Distinct 2 0.01%

Non-distinct 99998 99.99%

Values

Value	Frequency	Length	Percentage
1	81389	1	81.39%
2	18611	1	18.61%

Rule used: DQ_NAMOS_RIN_00001
Input columns: RRN, RRN_MARKETINGOVOJ_TRANZAKCI

Datatype

string(10) Documented

Decimal(1) 100.00%

Fixed Length String(1) 100.00%

Integer(1) 100.00%

String(1) 100.00%

Patterns

9 100.00%

Data Domain

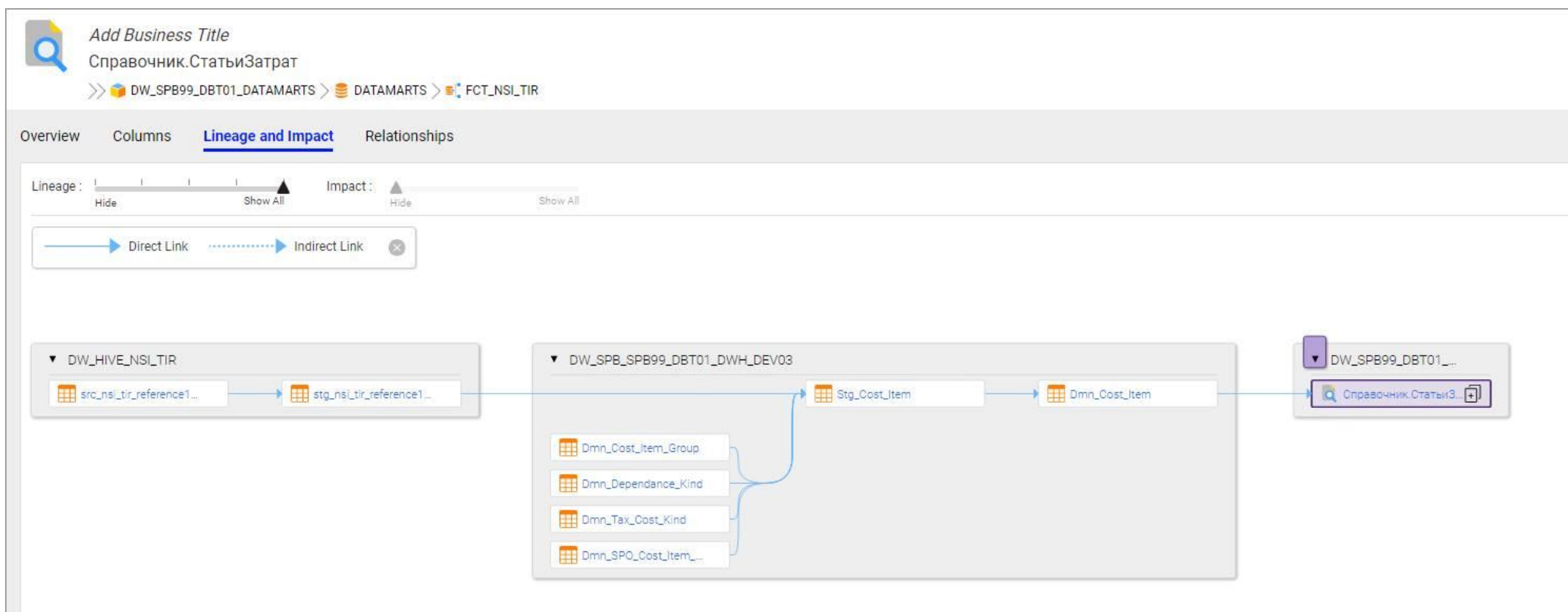
No data domain is inferred

Business Term

Data Preview

IDENFIKATO...	NOMER_SMENY	DATA_L_VREMY...	DATA_L_VREMY...	DATA_L_VREMY...	KOD_PRODUKTA	SUMMA	VID_VYBYTIYA	PLOTNOST	NOMER_BONU...	RRN	SUMMA_BEZ...	RRN_MARKETI...	DQ_Namos_PL...	DQ_Namos_VR...	DQ_Namos_VR...	DQ_Namos_NO...	DQ_NAMOS_S...	DQ_Namos_Su...	DQ_NAMOS_R...
1	'0000025061	920	10/07/2019 04.0...	10/08/2019 03.5...	10/07/2019 13.4...	4	1509.90000000	303	0.75360000		1509.95000000	0	1	1	1	0	1	1	0
2	'0000025061	920	10/07/2019 04.0...	10/08/2019 03.5...	10/07/2019 05.4...	99000000550060	95.00000000	0	0.00000000		95.00000000	0	1	1	1	0	1	1	0
3	'0000025061	920	10/07/2019 04.0...	10/08/2019 03.5...	10/07/2019 05.4...	3	1000.00000000	0	0.73000000		1000.13000000	0	1	1	1	0	1	1	0
4	'0000025061	920	10/07/2019 04.0...	10/08/2019 03.5...	10/07/2019 05.4...	990000000550300	106.00000000	0	0.00000000		106.00000000	0	1	1	1	0	1	1	0
5	'0000025061	920	10/07/2019 04.0...	10/08/2019 03.5...	10/07/2019 06.0...	13	1400.00000000	0	0.85223000		1400.24000000	0	1	1	1	0	1	1	0
6	'0000025061	920	10/07/2019 04.0...	10/08/2019 03.5...	10/07/2019 06.0...	990000000550061	122.00000000	0	0.00000000		122.00000000	0	1	1	1	0	1	1	0
7	'0000025061	920	10/07/2019 04.0...	10/08/2019 03.5...	10/07/2019 06.0...	9900000010409	90.00000000	0	0.00000000		90.00000000	0	1	1	1	0	1	1	0
8	'0000022206	901	10/07/2019 15.0...	10/08/2019 03.5...	10/07/2019 20.0...	9900000004296	162.00000000	0	0.00000000		162.00000000	0	1	1	1	0	1	1	0
9	'0000022206	901	10/07/2019 15.0...	10/08/2019 03.5...	10/07/2019 20.0...	99000000049542	118.00000000	0	0.00000000		118.00000000	0	1	1	1	0	1	1	0
10	'0000022206	901	10/07/2019 15.0...	10/08/2019 03.5...	10/07/2019 17.3...	4	600.00000000	0	0.74537000		600.28000000	0	1	1	1	0	1	1	0

Каталог данных подключается к компонентам хранения данных (Hadoop, MS SQL Server) и исполняемым ETL-приложениям, автоматически извлекая/обновляя цепочки происхождения данных и правила/формулы трансформации



Цепочки происхождения данных в каталоге разворачиваются до полей таблиц



Add Business Title

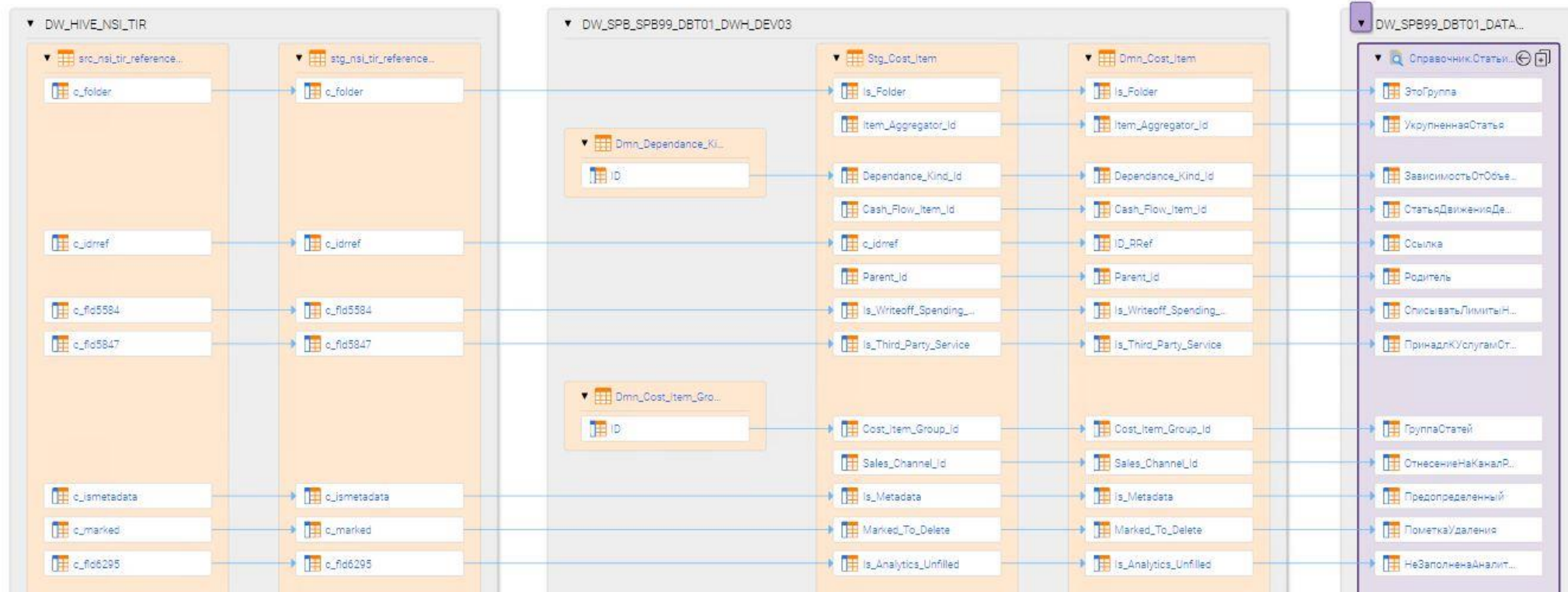
Справочник.СтатьиЗатрат

>> DW_SPB99_DBT01_DATAMARTS > DATAMARTS > FCT_NSL_TIR

Overview Columns **Lineage and Impact** Relationships

Lineage: Hide Show All Impact: Hide Show All

→ Direct Link Indirect Link ✕



Add Business Title
PRECALC_FULL
DW_SPB99_DBT02_DWH > DWH > MK

Overview Columns **Lineage and Impact** Relationships

PRECALC_FULL > app_mk_transaction_test

Direct Link Indirect Link

Транзакции_Розница_Цена_Реализ_Брутто_С_Налогом_Руб_Литр

Properties

RESOURCE NAME: app_mk_transaction_test
RESOURCE TYPE: BDMScanner
PRECISION: 18
INPUT: false
OUTPUT: true
SCALE: 7
UPPER: 1
GROUP NAME: Group
LOWER: 1
EXPRESSION: IIF(марНом2Скидка_Номенклатура_ЕКТУ_Не_Учитывать_Скидку = 'Да', b.Транзакции_Розница_Цена_Реализ_Нетто_С_Налогом_Руб_Литр, b.Транзакции_Розница_Цена_Реализ_Брутто_С_Налогом_Руб_Литр)
VARIABLE: false
LAST MODIFIED: 2019-10-01T10:09:023Z

Add Business Title
Справочник.СтатьиЗатрат
DW_SPB99_DBT01_DATAMARTS > DATAMARTS > FCT_NSL_TIR

Overview Columns **Lineage and Impact** Relationships

Справочник.СтатьиЗатрат > app_load_cost_item_full

Direct Link Indirect Link

Is_Affect_To_Tax_Base

Properties

RESOURCE NAME: app_load_cost_item_full
RESOURCE TYPE: BDMScanner
UPPER: 0
LOWER: 0
INPUT: false
OUTPUT: true
VARIABLE: false
PRECISION: 1
SCALE: 0
GROUP NAME: Exp_calc_hash_string_and_boolean_dft_vals
EXPRESSION: SUBSTR(SUBSTR(C_FLD5582_PASS, 2, 1), 0, 1)
LAST MODIFIED: 2019-06-11T15:38:46.355Z

Просмотр актуальных формул и правил трансформации на цепочках происхождения данных по каждому полю





Статья затрат

src_nsi_tir_reference1131

>> DW_HIVE_NSI_TIR > Hive Metastore > nsi_tir

Overview Columns Lineage and Impact Relationships

Name	Business Title	Data Domains	Null Distinct Non-Distinct %	Source Data Type Inferred Data Types
1 c_code	Код		0 100 0	string (N/A) Decimal(14) 0.13% +2 more
2 c_description	Наименование		0 24.66 75.34	string (N/A) String(149) 100.00%
3 c_fd2247rref	ВидРасходоВНУ		3.16 0.36 96.48	string (N/A) Decimal(32) 0.53% +3 more
4 c_fd2826rref	ВидСтатьиЗатратМСФОКР		3.16 0.81 96.03	string (N/A) Decimal(32) 3.86% +3 more
5 c_fd2827rref	ВидСтатьиЗатратМСФООАР		3.16 0.76 96.08	string (N/A) Decimal(32) 3.86% +3 more
6 c_fd2934rref	СтатьяДвиженияДенежныхСредств		3.16 5.43 91.41	string (N/A) Decimal(32) 85.42% +3 more
7 c_fd5426rref	ВидСтатьиЗатратМСФООР		3.16 0.64 96.20	string (N/A) Decimal(32) 4.05% +3 more
8 c_fd5580rref	ЗависимостьОтОбъемаПродаж		0 0.04 99.96	string (N/A) Decimal(32) 87.89% +3 more
9 c_fd5581rref	ВидУправления		0 0.04 99.96	string (N/A) Decimal(32) 87.12% +3 more
10 c_fd5582	ВлияниеНаВеличинуБазыНП		0 0.03 99.97	string (N/A) Decimal(2) 100.00% +3 more
11 c_fd5584	СписыватьЛимитыНаОсвоение		3.16 0.03 96.81	string (N/A) Decimal(2) 100.00% +3 more
12 c_fd5585rref	ВидТранспорта		3.16 0.11 96.73	string (N/A) Decimal(32) 99.28% +3 more
13 c_fd5586rref	ВидРынка		3.16 0.06 96.78	string (N/A) Decimal(32) 99.28% +3 more
14 c_fd5587	МожетБытьРБП		3.16 0.03 96.81	string (N/A) Decimal(2) 100.00% +3 more
15 c_fd5588	ИспользуетсяДляКонтроляЛимитовПоОУУУ		3.16 0.03 96.81	string (N/A) Decimal(2) 100.00% +3 more



Статья затрат

>> EDC_AXON

Overview Relationships

▼ Source Description

Группировки затрат, отражающих потребление производственных ресурсов по их видам, образующих себестоимость продукции (работ, услуг)

▼ Business Logic

Статьи затрат используются для понимания, на что были израсходованы средства предприятия, не обязательно только денежные, но и материалы, товарные или продуктовые запасы, износ оборудования и прочее.

▼ Related Glossary Assets (0) ⓘ

No details found.

▶ Classified Assets (0) ⓘ

▼ Related Technical Assets (4) ⓘ

Asset Name	Path
Stg_Cost_Item	DW_SPB_SPB99_DBT01_DWH_DEV03 / DWH_DEV03 / FCT_NSI_TIR
src_nsi_tir_reference1131	DW_HIVE_NSI_TIR / Hive Metastore / nsi_tir
Dmn_Cost_Item	DW_SPB_SPB99_DBT01_DWH_DEV03 / DWH_DEV03 / FCT_NSI_TIR
stg_nsi_tir_reference1131	DW_HIVE_NSI_TIR / Hive Metastore / nsi_tir

Профилирование позволяет выявить:

- типы данных
- маски
- значения и с какой частотой они встречаются в данном поле
- степень схожести с полями других таблиц по следующим параметрам: наименования, похожесть значений, похожесть типов, похожесть масок

Код
c_code

Overview Lineage and Impact Relationships

Description
Код статьи затрат в системе НСИ-ТИР ДРП

Value Frequency
Total Rows: 6037
Max: -30.1.2
Min: 00

Value	Frequency	Percentage
00	1	0.02
01	1	0.02
02000000000000	1	0.02
03000000000000	1	0.02
04000000000000	1	0.02
05000000000000	1	0.02
05000000000001	1	0.02
05000000000002	1	0.02

Similar Columns

Column Name	Confidence
.../stg_nsi_tir_reference1131 c_code	99%
.../Dmn_Cost_Item Код Статьи затрат НСИ-ТИР Code	99%
.../Vw_Dmn_Cost_Item Code	93%

Pattern

Pattern	Rows	Percentage
99,999,999,9(4),999	5326 rows	88.22%
Others	711 rows	11.78%

Inferred Data Types

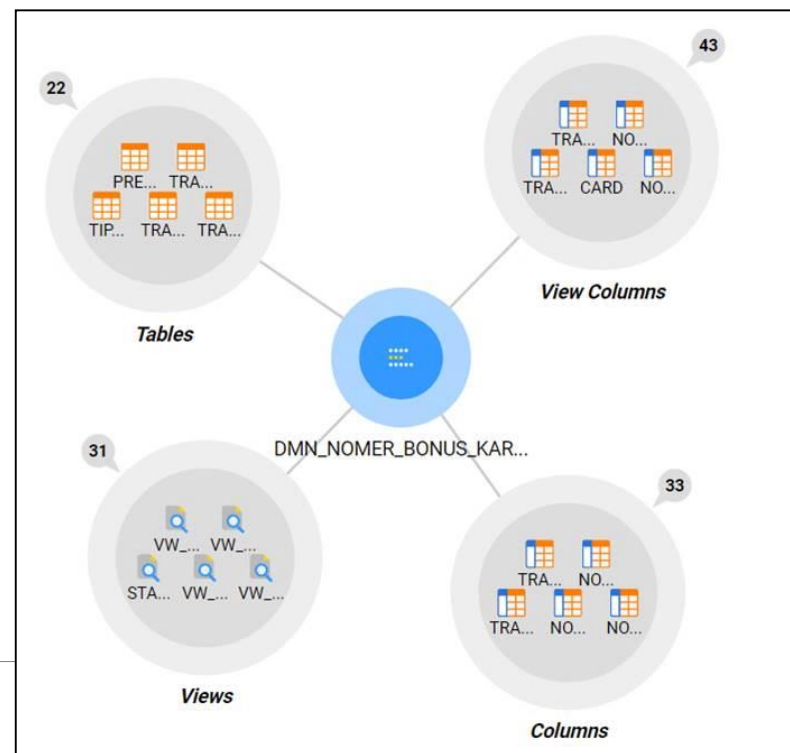
Data Type	Rows	Percentage
Decimal(14)	8 rows	0.13%
Integer(14)	8 rows	0.13%
String(19)	6037 rows	100.00%

Применение разметки данных: автоматическое определение доменов, принадлежности к коммерческой тайне и персональным данным

DMN_NOMER_BONUS_KARTY

Проверяет данные по маскам бонусных карт

Resource Name: **DataDomain**



Add Business Title
PDN

Overview Lineage and Impact Relationships

Description

Композитный домен проверяет таблицы на вхождение персональных данных о клиенте: ФИО, дата рождения, номер телефона, E-mail

Source Description

Композитный домен с содержанием пунктов КТ-040 (ПДН)

Found in

Name	Resource Type	Business Term
DW_SPB99_DBT02_DWH	JDBC	
Instance 1 :		
VLADELCOY_BONUSNYX_KART_OLD		

Вид выбытия ТЕРМИН
Glossary

SUMMARY RELATIONSHIPS DATA QUALITY STAKEHOLDERS IMPACT DATA ENTERPRISE CATALOG HISTORY CHANGE

DEFINITION

Вид выбытия нефтепродуктов и сопутствующих товаров и услуг с АЗС

Ref.: GLOS-130

Alias Names: Not specified

Format: Другое - Not specified

LDM Reference: Not specified

Business Logic: Вид выбытия используется для выделения оборотов в коммерческом учете. Источниками транзакционных данных являются различные АСУ АЗС: Namos, Мобильная карта (МК) и др. Данные, приходящие из этих систем имеют различные форматы и значения характеристики "Вид выбытия". Для выделения общих групп оборотов в коммерческом учете производится приведение данных к стандартизированному значению, которые хранятся в системе НСИ-ТИР ДРП.

Алгоритм определения вида выбытия для данных СТ и НП из АСУ АЗС Namos:

Мегтинг видов выбытия в транзакционных данных НП АСУ АЗС МК

Вид выбытия / Именованное	Вид выбытия ID ТИР
На личный расчет (н/продукты)	A467343806AADA30433E4A93486773A2
На личный расчет со скидкой	9A8594C2855281644807A44A08A070DF
За бонусы (НП)	98230728731337A94F6F5A86F5088C1B
На личный расчет со скидкой	9A8594C2855281644807A44A08A070DF
На личный расчет по бонусным картам	888E098580735F1C43035098D78D03EF
Снятие ПЛ	A467343806AADA30433E4A93486773A2
На личный расчет (н/продукты)	A467343806AADA30433E4A93486773A2
Бонусная карта по цене стелы	A3D2A980CB0134104138D02D72C50218
Бонусная карта со скидкой	9469042445638A5F41E797C4A37E5D3A

- ведение паспортов бизнес-сущностей, определений, методик расчета
- привязка к владельцам, стюардам и т.п.
- привязка к бизнес-процессам, тегам
- ведение истории изменений

Карта ПЛ ТЕРМИН
Glossary

SUMMARY RELATIONSHIPS DATA QUALITY STAKEHOLDERS IMPACT DATA ENTERPRISE CATALOG HISTORY CHANGE DIFFERENCE REPORT x

HISTORY

From Date: 01-Apr-2019 To Date: 08/10/2019 Compare

Object	Update Type	Field	From	To	Author	Date
Glossary	Updated	Description	Карта ПЛ "Нам по пути": Карта, выдаваемая Участнику ПЛ, которая пред	Карта ПЛ "Нам по пути": Карта, выдаваемая Участнику ПЛ, которая пред		12-Apr-2019
Glossary	Status Change	Lifecycle	На согласовании	Согласован		12-Apr-2019
Glossary	Updated	Description	Карта ПЛ "Нам по пути": Карта, выдаваемая Участнику ПЛ, которая пред	Карта ПЛ "Нам по пути": Карта, выдаваемая Участнику ПЛ, которая пред		10-Apr-2019
Stakeholders	Updated	reltype	39			12-Apr-2019
Stakeholders	Added	Name				12-Apr-2019
Stakeholders	Added	Role		Стюард объекта гlossария		12-Apr-2019
Stakeholders	Status Change	Accepted				12-Apr-2019
Stakeholders	Added	Axon Status		Active		12-Apr-2019

8 records



DS-166: fct_regional_index

Data Set in DW_SPB_SPB99_DBT01_DWH_DEV03

Edit

- SUMMARY
- ATTRIBUTES
- RELATIONSHIPS
- STAKEHOLDERS
- VALUES
- DATA QUALITY
- IMPACT
- HISTORY
- CHANGE

FOLLOW

MAP

Map type: Data Set Lineage | Layout: Left-To-Right | Overlay: Data Quality | Filters: All selected (2)



INBOUND RELATIONSHIPS

Attribute Name	Source System	Related Data Set	Related Attribute	Type	Scope
index_value	DW_Hive_internet_sources	DS-163: stg_regional_index_pqt	index_value	Источник	
dt	DW_Hive_internet_sources	DS-163: stg_regional_index_pqt	index_date	Источник	



DS-159: fct_stock_index

Data Set in DW_SPB_SPB99_DBT01_DWH_DEV03

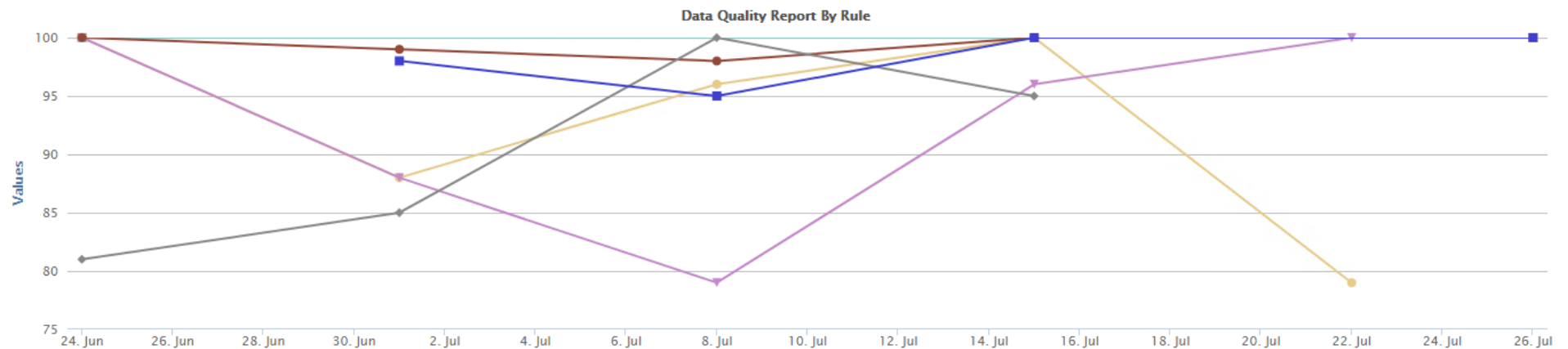
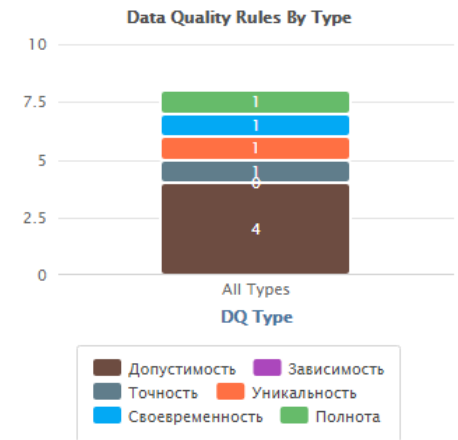
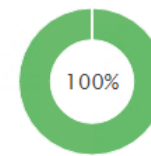
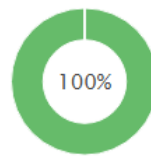
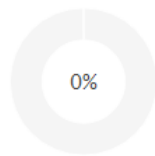
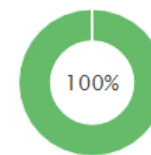
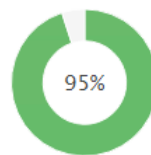
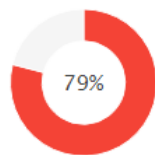
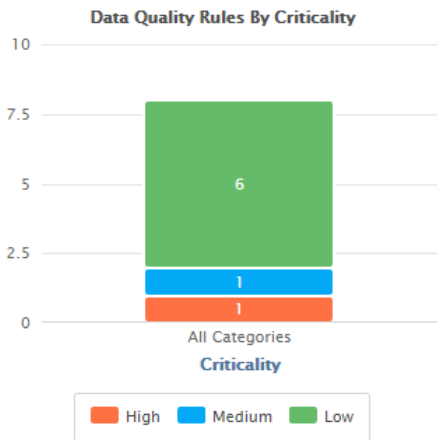
Edit

SUMMARY ATTRIBUTES RELATIONSHIPS STAKEHOLDERS VALUES **DATA QUALITY** IMPACT HISTORY CHANGE

FOLLOW

Dashboard Rules

DATA QUALITY

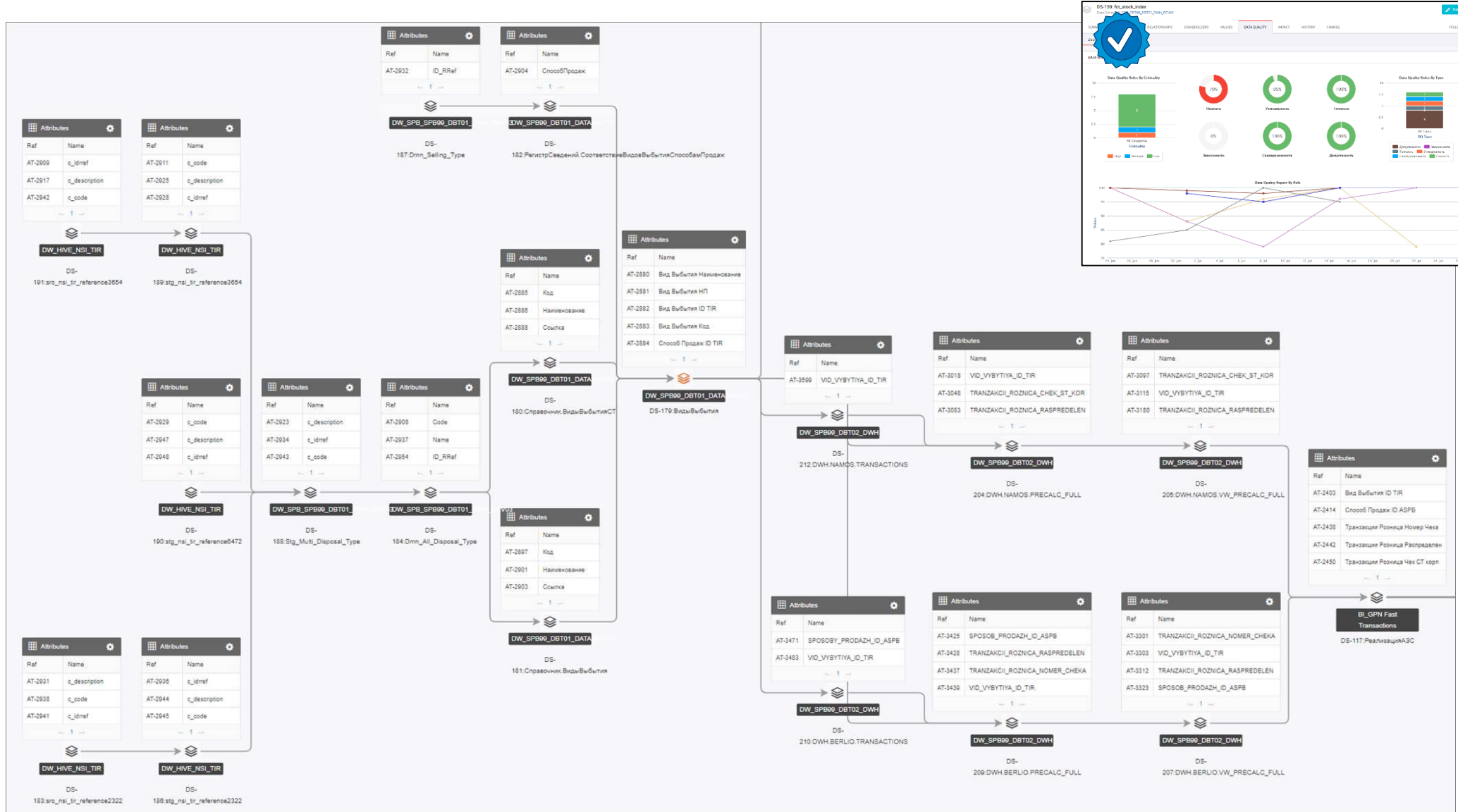


До мая 2019 года разработка «озера данных» велась без использования ETL-средств Informatica (ручной код, MS SSIS)

Для этого сценария функционал системы управления данными работает в режиме ручного заполнения

Сравнение функциональности управления данными для сценариев:

Функционал	Целевой сценарий: ETL Informatica	Сценарий: ETL другими средствами
Сбор и актуализация метаданных источников, ведение единого каталога данных	ДА в EDC в автоматическом режиме	ДА в Ахон в ручном режиме
Ведение описаний в карточках БД/таблиц/полей	ДА в EDC в ручном режиме	ДА в Ахон в ручном режиме
Построение Lineage на уровне таблиц/полей	ДА в EDC в автоматическом режиме	ДА в Ахон в ручном режиме
Просмотр формул, скриптов трансформации данных	ДА в EDC	НЕТ
Привязка технических метаданных к бизнес-терминам	ДА в EDC и Ахон в ручном режиме	ДА в Ахон в ручном режиме
Автоматическое профилирование данных: уникальность, примеры значений, похожесть, типы данных, маски данных	ДА в EDC в автоматическом режиме	НЕТ
Автоматическая разметка на КТ, ПДн	ДА в EDC в автоматическом режиме	ДА в Ахон в ручном режиме



Технически работают подключения к источникам с построением статусов, интегральных дашбордов, динамики правил по качеству данных

1. **Зона сырых данных** - сырой слой (копии данных из источников), где данные имеют имена как в источниках

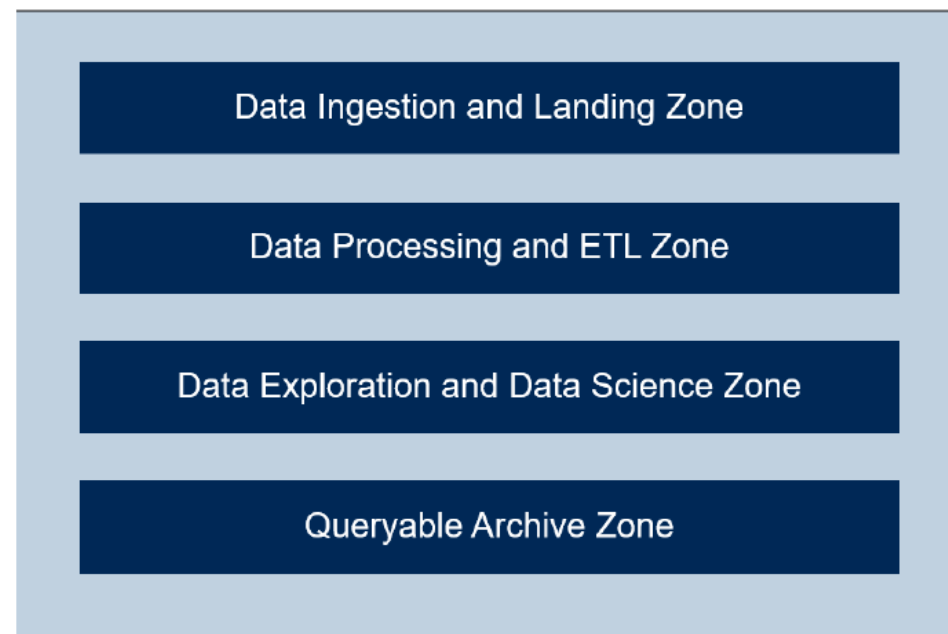
2. **Продуктивная фабрика данных** - слой продуктивных трансформаций для аналитического хранилища с фокусом на оптимизацию, производительность и управление, объекты данных имеют целевые наименования и паспортизируются в Системе управления данными

3. **Исследовательская лаборатория** - зона «быстрой разработки» для проведения исследований и формирования прототипов решений, фокус на скорость загрузки и трансформации данных, минимум формализации, выбор инструментов не ограничивается, наименования объектов данных не регулируются, разработка может проводиться сотрудниками-аналитиками

Контур управления данными

Лучшая практика (Gartner)*

Data Lake Zones



ID: 367848

© 2018 Gartner, Inc.

Управление данными

- ✓ Автоматически обновляемые цепочки происхождения данных от источника до пользовательских витрин и BI-приложений
- ✓ Бизнес-гlossарий, интегрированный с уровнем представлений хранилища данных (поля отчетов и моделей расшифровываются)
- ✓ Единый управляемый каталог правил и проверок по качеству данных
- ✓ Бесшовная интеграция правил по качеству данных в ETL-цепочки
- ✓ Пользовательский портал для просмотра glossария, происхождения данных и правил по качеству данных

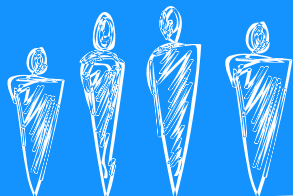
Озеро данных

- ✓ Снижение затрат на интеграцию данных для всех связанных с аналитикой проектов (KPI – не менее 50% всех проектов ДРП на конец 2019 года используют инфраструктуру Озера данных)
- ✓ Доступность данных на всех слоях конвейера для пользователей
- ✓ Интеграция с инструментами Data Science (Лаборатория + Spark + In-database Python, R)
- ✓ Промышленные инструменты ETL и хранения для всех типов данных, потоковая и пакетная обработка
- ✓ Промышленные инструменты для создания интерфейсов интеграции

Задачи



Интеграция компонент
Интеграция с DevOps
Производительность и стабильность



18 сотрудников в команде, без внешних подрядчиков

Хранилище данных (январь - сентябрь 2019)

10 T6 объем данных
1350 таблиц и view
2000 Заданий в сутки

Фокусы



Только актуальные бизнес-задачи
Управляем только важными данными
Стандартизация практик и процессов



Период	Название подразделения	Орг. Изменение	Численность сотрудников
2012 – 2013 гг	Единое ответственное лицо за проекты BI и HSI	[все работы выполняются внешними подрядчиками]	1
2014 год	Группа развития BI	Новые роли: <ul style="list-style-type: none"> • «Архитектор BI» • «Разработчик BI» • «Системный аналитик BI» 	4
2015 – 2016 годы		Новая роль <ul style="list-style-type: none"> • «Эксперт по качеству данных» • «Администратор BI» 	10
2017 год	Центр компетенций BI	[Рост внутреннего ресурса, отказ от внешних подрядчиков] Новая роль <ul style="list-style-type: none"> • «Архитектор данных» 	30
2018 год	Офис управления данными	[Старт создания системы управления данными, сразу ставка на внутренние ресурсы] Новые роли: <ul style="list-style-type: none"> • «Эксперт по управлению метаданными» • «Математик-программист» (Data scientist) • «Системный аналитик Озера данных» • «Архитектор Озера данных» 	40
2019 год	Центр аналитических решений	+ Центр компетенций по Data Science + Центр развития компетенций в аналитике Новые роли: <ul style="list-style-type: none"> • «Тренер по аналитике» • «Разработчик Озера данных» • «Администратор систем управления данными» 	60

1. Разработка аналитических решений

BI-приложения
Хранилища данных для аналитики
Решения для «больших данных»

6. Обучения и тесты по методам и инструментам анализа данных

Python, R, SQL, BI
Методы анализа данных
Машинное обучение и т.п.

5. Управление централизованной аналитической инфраструктурой

Аналитические песочницы и инструменты
Инжиниринг и окружение
Стандартизация: данные как сервис

2. Управление данными

Управление качеством данных
Единый каталог данных и бизнес-словарь
База знаний и консультирование по данным

3. Центр компетенций Data Science

Модели и прототипы на данных
Участие в проектах с аналитической составляющей
Центр знаний и ресурсов

4. Бизнес-партнерство в части новых проектов

Оценка новых инициатив с точки зрения данных
Участие в проектных командах





Иван Черницын

CHERNITSYN.IG@gazprom-neft.ru

Руководитель Центра аналитических решений
Дирекция региональных продаж
ПАО «Газпром нефть»

