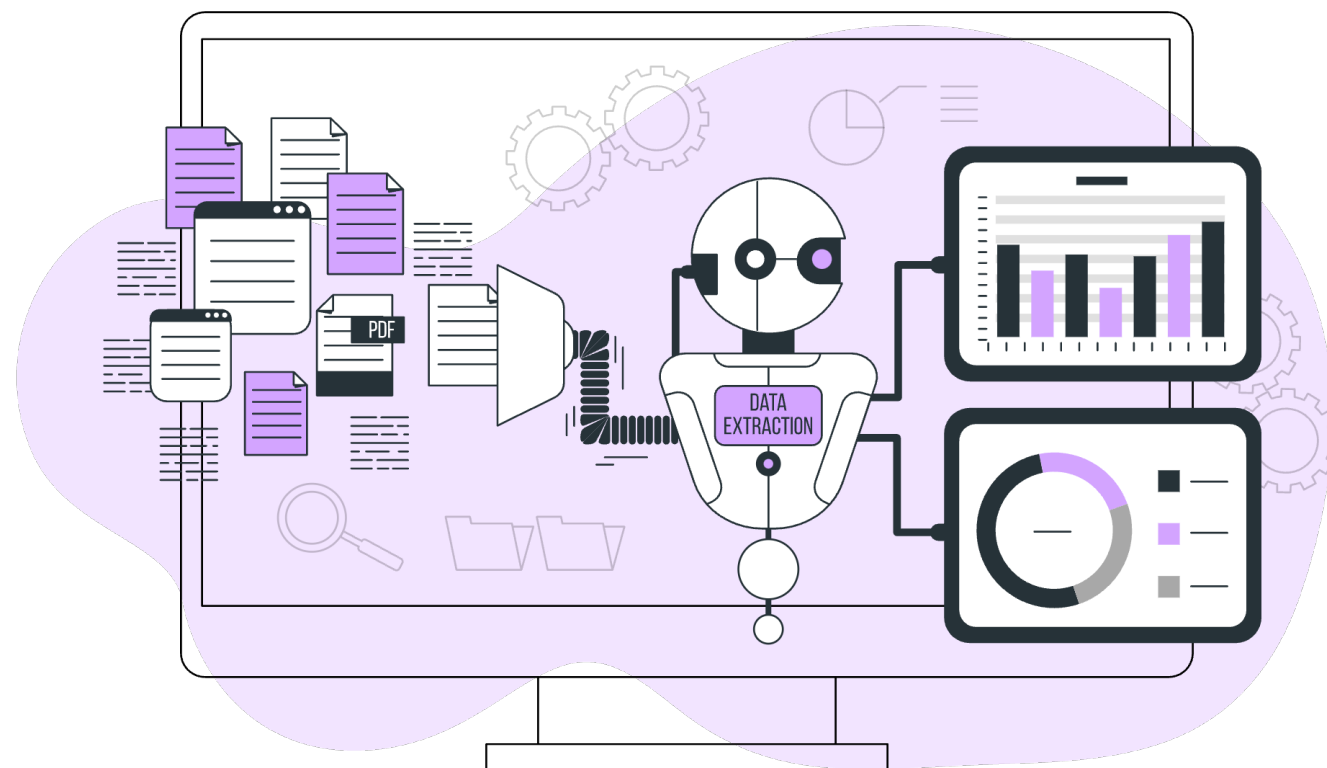


# Управление данными и разработкой как основа для применения ИИ

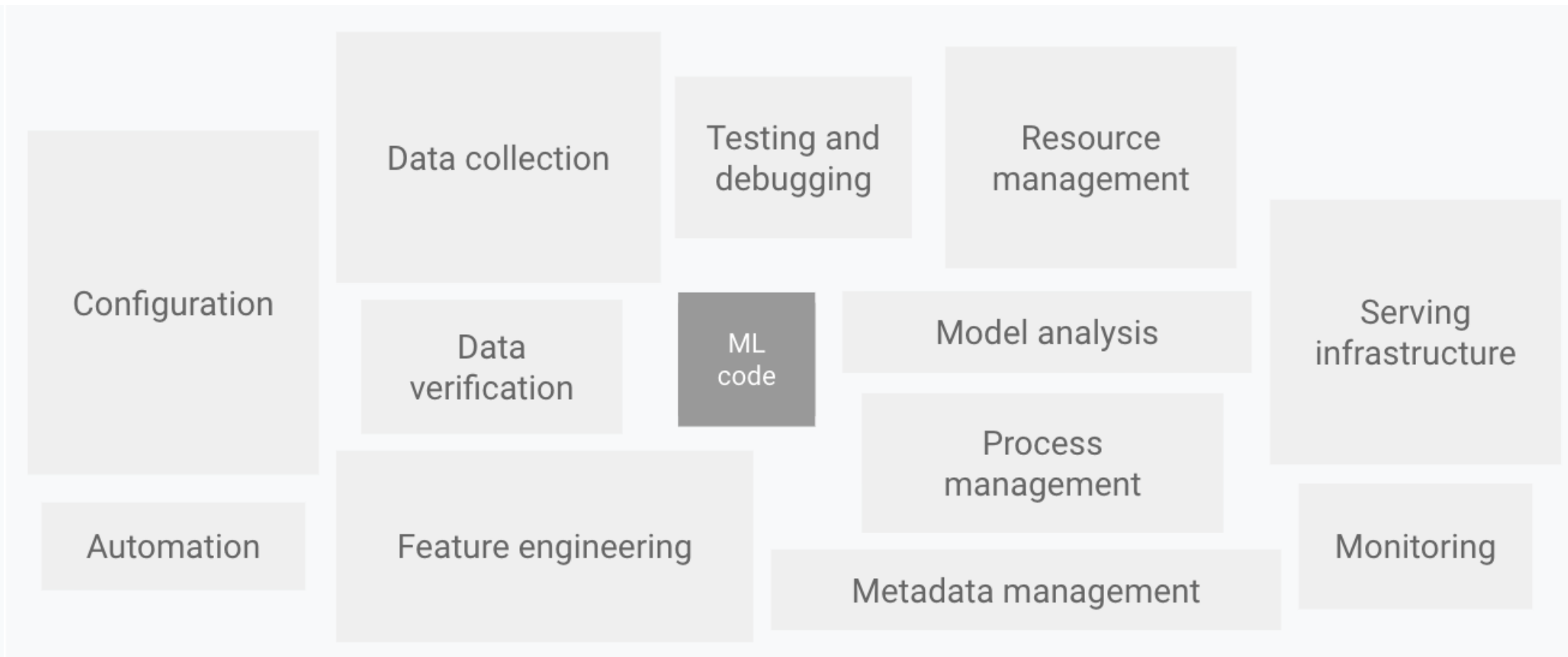


**Загорулькин Дмитрий Эдуардович**

Заместитель директора Центра стратегической  
аналитики и больших данных  
[dzagorulkin@hse.ru](mailto:dzagorulkin@hse.ru)



# Составные части ML системы



[https://proceedings.neurips.cc/paper\\_files/paper/2015/file/86df7dcfd896fcaf2674f757a2463eba-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2015/file/86df7dcfd896fcaf2674f757a2463eba-Paper.pdf)

# Создание AI проектов

01

Процессы

02

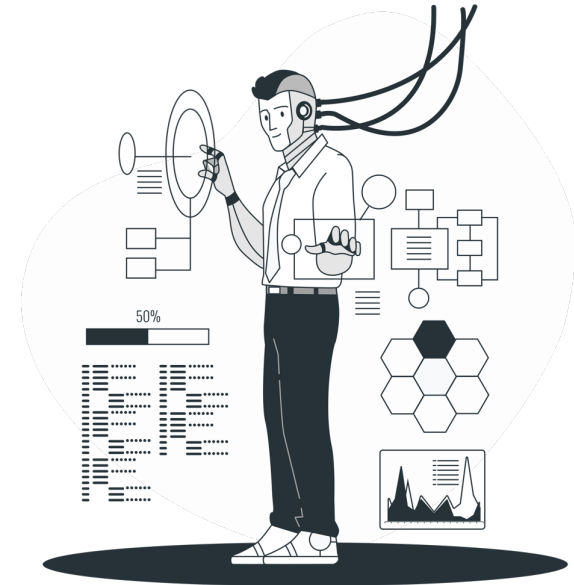
Инфраструктура

03

Данные

04

Моделирование



Fall into ML  
'24

# Понятная постановка задачи

Нам нужно сделать модель, которая определяет уровень технологии.



ПЛОХАЯ ПОСТАНОВКА ЗАДАЧИ



Мы собираем данные об упоминании различных технологий в СМИ. Данные собираются регулярно. Нам необходимо улучшить точность оценки минимум на 20%. Также нужно учитывать, что технологии могут относиться к разным отраслям и иметь разный жизненный цикл. Уровни нужно агрегировать с 10 до 3 классов.

ХОРОШАЯ ПОСТАНОВКА ЗАДАЧИ



# Метрики — они для всех разные!



## 01

Бизнес сфокусирован на получении прибыли!  
Бизнес оперирует понятными метриками, такими как процент оттока клиентов, время, проведенное в продукте, количество пользовательских входов в месяц и т.д.

## 03

Для получения профита от внедрения ML решения важно связать модельные метрики и метрики бизнеса  
Важно найти и протестировать это влияние как можно скорее

## 02

Дата-сайентист оперирует модельными метриками (P/R, Accuracy, F1 и другие)

Высокие модельные метрики не гарантируют улучшение бизнес-метрик!

# Управление AI проектом



## SCRUM

Подходит под ИТ-проекты

## CRISP-DM (Cross-Industry Standard Process for Data Mining)

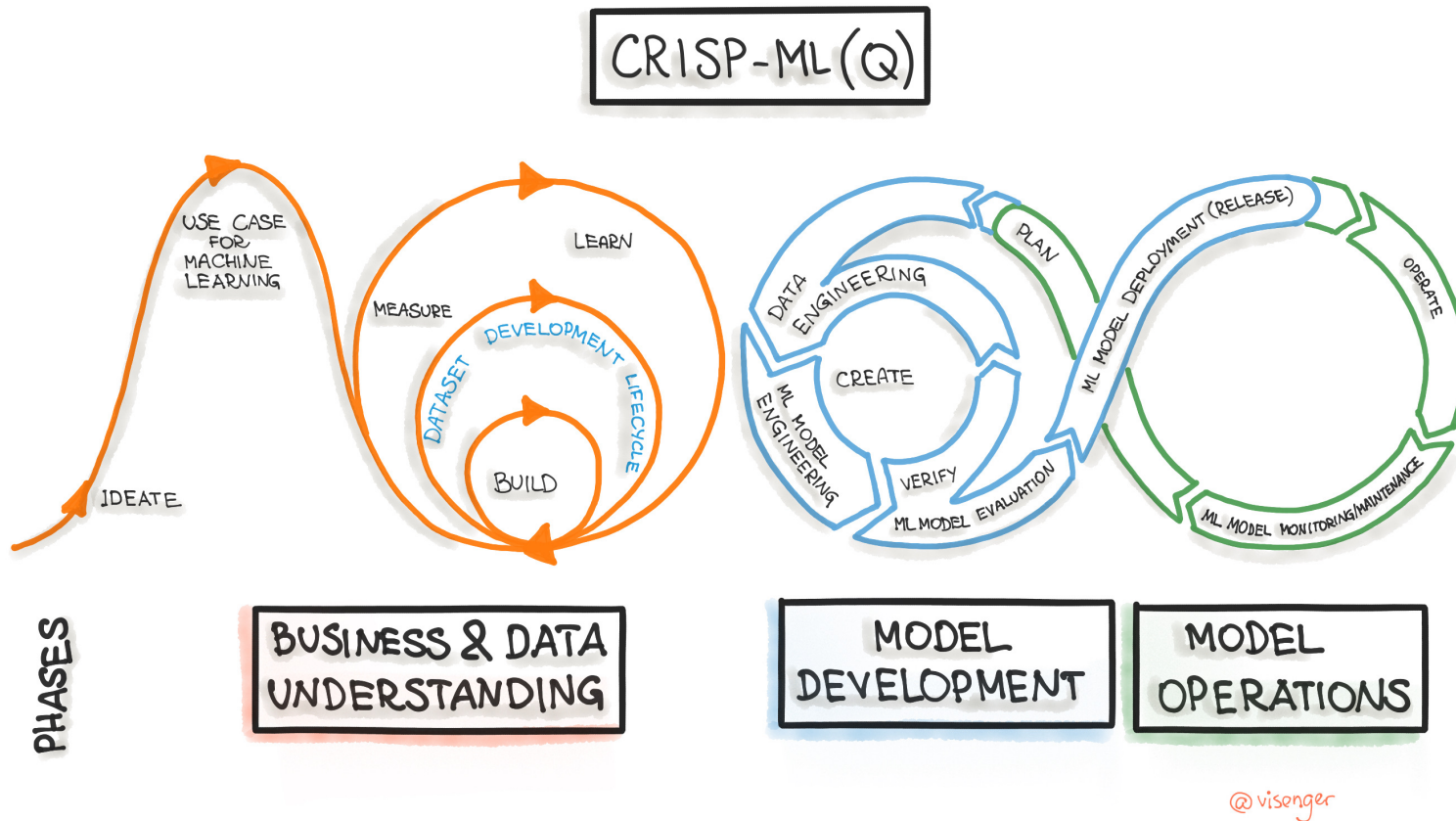
1. Понимание бизнес-целей (Business Understanding) – необходимо привлечение всех заинтересованных сторон
2. Понимание данных (Data Understanding) – проведение разведочного анализа и другие проверки данных
3. Подготовка данных (Data Preparation) – консолидация, агрегация
4. Моделирование (Modeling) – выбор методологии и построения модели
5. Оценка результата (Evaluation)
6. Внедрение модели/процесса (Deployment)

**Проблемы:** Не предназначена под итеративный процесс улучшения. Отсутствует QA.

## CRISP-ML(Q)

Расширение методологии

# CRISP-DM сравнение с CRISP-ML(Q)



CRISP-ML(Q)	CRISP-DM
Business & Data Understanding	Business Understanding
Data Preparation	Data Understanding
Modeling	Data Preparation
Evaluation	Modeling
Deployment	Evaluation
Monitoring & Maintenance	Deployment
	-

<https://ml-ops.org/content/crisp-ml> | <https://arxiv.org/pdf/2003.05155>

# Данные



Основа построения систем машинного обучения

Существуют разные типы данных. Для одних моделей нужны сильно структурированные данные, для других — нет (картинки, видео, книги и таблицы в БД)

Отсутствие культуры работы с данными зачастую является тормозом внедрения систем машинного обучения в бизнес-процессы

Необходимо развивать культуру работы с данными в части автоматической проверки качества данных и других data governance подходов (ведение дата-каталогов и управление метаданными и др.)

Не все данные подходят для моделирования!



# Управление данными



## DWH (Data Warehouse)

Сложность внесения изменений • Долго строить • Нужно хорошо понимать данные • В основном для структурированных данных

## DataLake

Единая точка получения данных • Нужно следить за качеством и метой • Нет ACID транзакций, эволюции схем

**DataLakehouse = DataLake + DWH**

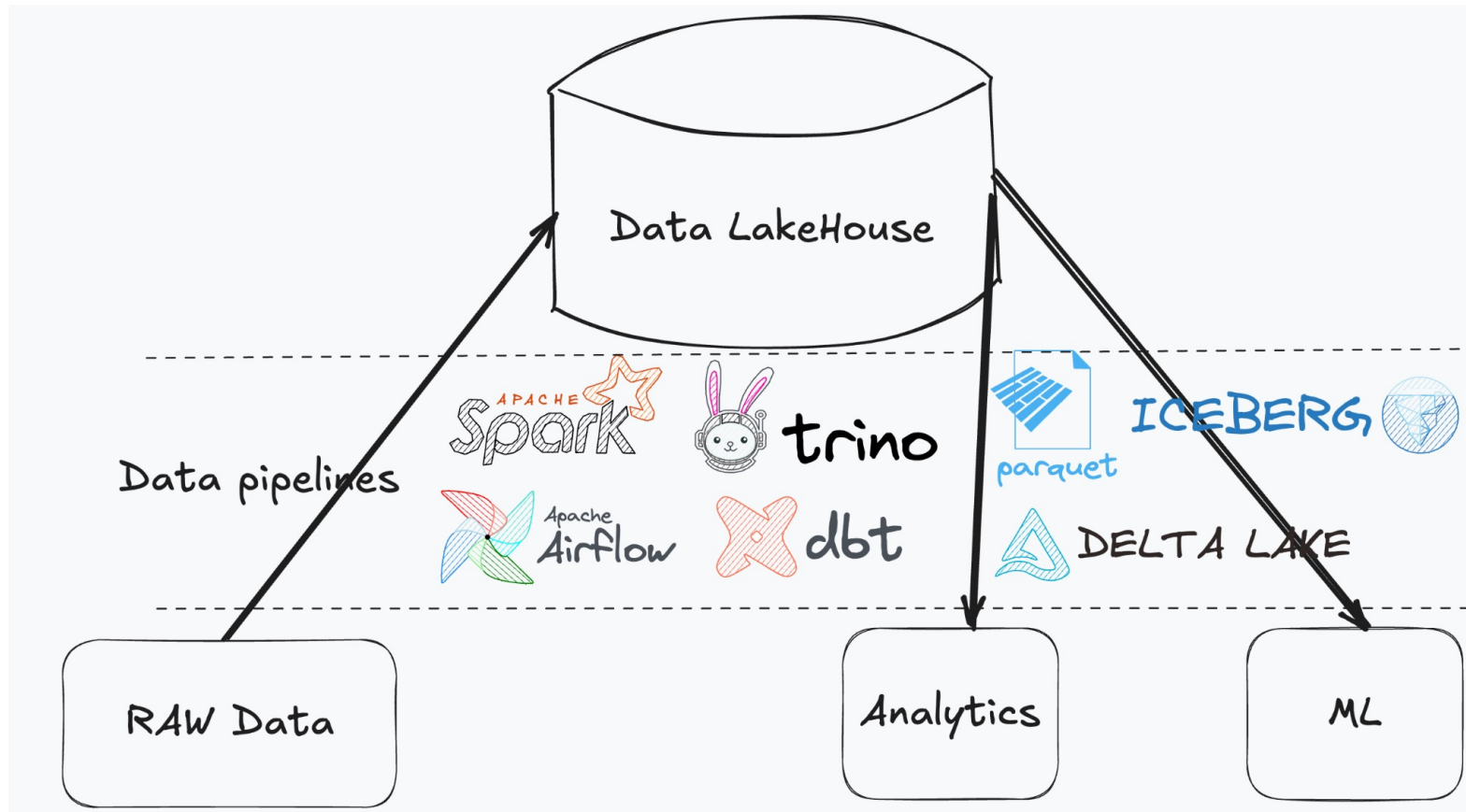
## DataMesh

Организационный подход к управлению данными • Полностью децентрализованный • За отдельные данные отвечает владелец, который передает их при необходимости потребителям

## DataFabric

Включает в себя элементы DataLakehouse и DataMesh и вносит дополнительно элементы data governance

# Процесс работы с данными



Fall into ML  
'24

# Инфраструктура для ML (DL)



Критерий	Своя	Облачная
Обслуживание инфраструктуры	Нужны квалифицированные инженеры (высокие операционные издержки)	–
Инженеры DevOps/MLOps	+	+
Готовность крупного бизнеса передавать чувствительные данные (безопасность)	Нужны специалисты по безопасности	Нет
Расширяемость (Managed решения)	- / +	Ограничена предоставляемыми сервисами
Vendor lock	Нет	Да
Гибкая масштабируемость	Ограничена физическими серверами и используемыми инструментами	Да
Отказоустойчивость	Нужно несколько ЦОДов	+
Цена (особенно при использовании GPU ускорителей)	Получается на порядок дешевле на длинной дистанции	Pay as you go

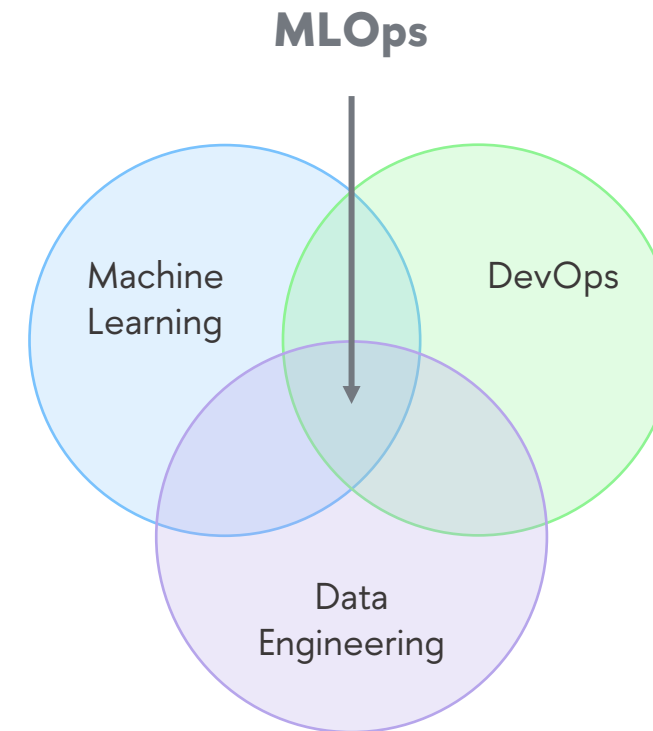
**Вывод:** использование только облачной инфраструктуры для DL выходит дорого, идеально использовать гибрид, если это возможно для бизнеса. Разработка внутри, инференс для клиентов снаружи.

**Всегда проводите оценку, какие ресурсы потребуются под вашу систему и будет ли готов бизнес на такие расходы!**

## Основные плюшки:

- Автоматизация рутинных процессов
- Масштабируемость
- CI/CD для моделей (canary deploy)
- Воспроизводимость
- Улучшения взаимодействия между командами (DS, OPS, DEV)
- Мониторинг решений

Большое количество инструментов разного уровня зрелости.  
Может быть сложно интегрировать их между собой и в бизнес-процессы компании.



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

### INFRASTRUCTURE

**STORAGE**  
Amazon S3, Google Cloud Storage, IBM Cloud Storage, Pure Storage, Veeva, NetScout Systems, Cohesity, Vast, DataEa, Qumulo

**HADOOP**  
Cloudera, Amazon EMR, Databricks, AWS EMR, Hadoop Distribution, Pivotal, Jethro, Cloud Platform, IBM Watson Analytics, QRM, IBM InfoSphere

**DATA LAKES**  
Databricks, Amazon Data Lake Storage, AWS Lake Formation, AWS Power Systems, IBM Watson Analytics, QRM, IBM InfoSphere

**DATA WAREHOUSES**  
Amazon Redshift, Snowflake, Google BigQuery, Microsoft Azure Synapse Analytics, Oracle Exalytics, SAP HANA, IBM Watson Analytics, QRM, IBM InfoSphere

**STREAMING / IN-MEMORY**  
Amazon Kinesis, Databricks, Apache Spark, Amazon EMR, Google Cloud Dataflow, Microsoft Azure Databricks, IBM Watson Analytics, QRM, IBM InfoSphere

**NOSQL DATABASES**  
Google Cloud Spanner, Amazon DynamoDB, Amazon DocumentDB, SAP HANA, Oracle NoSQL Database, MongoDB, Oracle, Couchbase, IBM Watson Analytics, QRM, IBM InfoSphere

**REAL TIME DATABASES**  
Amazon Neptune, IBM Watson Analytics, QRM, IBM InfoSphere

**GRAPH DBs**  
Amazon Neptune, IBM Watson Analytics, QRM, IBM InfoSphere

**RDDBMS**  
Oracle, SAP HANA, IBM Watson Analytics, QRM, IBM InfoSphere

**MPP DBs**  
Teradata, Vertica, IBM Watson Analytics, QRM, IBM InfoSphere

**ETL / ELT / DATA TRANSFORMATION**  
dbt, Talend, Alteryx, Informatica, SAP HANA, IBM Watson Analytics, QRM, IBM InfoSphere

**REVERSE ETL**  
Census, IBM Watson Analytics, QRM, IBM InfoSphere

**DATA INTEGRATION**  
MuleSoft, Tealium, Inoplogic, IBM Watson Analytics, QRM, IBM InfoSphere

**DATA GOVERNANCE & ACCESS**  
IBM Watson Analytics, QRM, IBM InfoSphere

**PRIVACY & SECURITY**  
Very Good Security, Privabera, IBM Watson Analytics, QRM, IBM InfoSphere

**DATA OBSERVABILITY**  
Datastax, Monte Carlo, IBM Watson Analytics, QRM, IBM InfoSphere

**MGMT / MONITORING**  
Amazon CloudWatch, New Relic, Splunk, IBM Watson Analytics, QRM, IBM InfoSphere

**SERVERLESS**  
Amazon Lambda, IBM Watson Analytics, QRM, IBM InfoSphere

**CLUSTER SVCS**  
Amazon EC2, IBM Watson Analytics, QRM, IBM InfoSphere

### ANALYTICS

**BI PLATFORMS**  
Looker, Tableau, Power BI, SAP, IBM Watson Analytics, QRM, IBM InfoSphere

**VISUALIZATION**  
Tableau, Power BI, SAP, IBM Watson Analytics, QRM, IBM InfoSphere

**DATA ANALYST PLATFORMS**  
Microsoft, Pentaho, Alteryx, IBM Watson Analytics, QRM, IBM InfoSphere

**AUGMENTED ANALYTICS**  
ThoughtSpot, IBM Watson Analytics, QRM, IBM InfoSphere

**DATA CATALOG AND DISCOVERY**  
Metaphor, Atlan, IBM Watson Analytics, QRM, IBM InfoSphere

**LOG ANALYTICS**  
Splunk, IBM Watson Analytics, QRM, IBM InfoSphere

**QUERY ENGINE**  
Dremio, IBM Watson Analytics, QRM, IBM InfoSphere

**SEARCH**  
Elasticsearch, Oracle, IBM Watson Analytics, QRM, IBM InfoSphere

### MACHINE LEARNING & ARTIFICIAL INTELLIGENCE

**DATA SCIENCE NOTEBOOKS**  
JupyterLab, Binder, Colab, IBM Watson Analytics, QRM, IBM InfoSphere

**DATA SCIENCE PLATFORMS**  
DataRobot, Dataiku, IBM Watson Analytics, QRM, IBM InfoSphere

**ML PLATFORMS**  
DataRobot, Dataiku, IBM Watson Analytics, QRM, IBM InfoSphere

**DATA GENERATION & LABELLING**  
Amazon Mechanical Turk, Hive, IBM Watson Analytics, QRM, IBM InfoSphere

**MODEL BUILDING**  
Weights & Biases, IBM Watson Analytics, QRM, IBM InfoSphere

**FEATURE STORE**  
Feast, IBM Watson Analytics, QRM, IBM InfoSphere

**DEPLOYMENT & PRODUCTION**  
DataRobot, Dataiku, IBM Watson Analytics, QRM, IBM InfoSphere

**MODEL MONITORING & OBSERVABILITY**  
Arthor, IBM Watson Analytics, QRM, IBM InfoSphere

**COMPUTER VISION**  
Microsoft Azure, IBM Watson Analytics, QRM, IBM InfoSphere

**SPEECH**  
Siri, Amazon Alexa, IBM Watson Analytics, QRM, IBM InfoSphere

**NLP**  
IBM Watson Analytics, QRM, IBM InfoSphere

**SYNTHETIC MEDIA**  
DeepBrain AI, IBM Watson Analytics, QRM, IBM InfoSphere

**HORIZONTAL AI**  
IBM Watson Analytics, QRM, IBM InfoSphere

**GPU DBS & CLOUD**  
Kinetic, IBM Watson Analytics, QRM, IBM InfoSphere

**AI HARDWARE**  
Google TPU, ARM, IBM Watson Analytics, QRM, IBM InfoSphere

### APPLICATIONS - ENTERPRISE

**SALES**  
Salesforce, IBM Watson Analytics, QRM, IBM InfoSphere

**MARKETING B2B**  
App Annie, IBM Watson Analytics, QRM, IBM InfoSphere

**MARKETING - B2C**  
Google Analytics, Tealium, ActionIQ, IBM Watson Analytics, QRM, IBM InfoSphere

**CUSTOMER EXPERIENCE / SERVICE**  
Qualtrics, SurveyMonkey, IBM Watson Analytics, QRM, IBM InfoSphere

**HUMAN CAPITAL**  
IBM Watson Analytics, QRM, IBM InfoSphere

**LEGAL**  
Ravel, IBM Watson Analytics, QRM, IBM InfoSphere

**REGTECH & COMPLIANCE**  
IBM Watson Analytics, QRM, IBM InfoSphere

**FINANCE**  
IBM Watson Analytics, QRM, IBM InfoSphere

**AUTOMATION & RPA**  
IBM Watson Analytics, QRM, IBM InfoSphere

**SECURITY**  
IBM Watson Analytics, QRM, IBM InfoSphere

**ADVERTISING**  
Xandr MediaMath, IBM Watson Analytics, QRM, IBM InfoSphere

**EDUCATION**  
IBM Watson Analytics, QRM, IBM InfoSphere

**REAL ESTATE**  
Redfin, IBM Watson Analytics, QRM, IBM InfoSphere

**GOVT & INTELLIGENCE**  
IBM Watson Analytics, QRM, IBM InfoSphere

**COMMERCE**  
IBM Watson Analytics, QRM, IBM InfoSphere

**FINANCE - LENDING**  
IBM Watson Analytics, QRM, IBM InfoSphere

**INSURANCE**  
IBM Watson Analytics, QRM, IBM InfoSphere

**HEALTHCARE**  
IBM Watson Analytics, QRM, IBM InfoSphere

**LIFE SCIENCES**  
IBM Watson Analytics, QRM, IBM InfoSphere

**TRANSPORTATION**  
IBM Watson Analytics, QRM, IBM InfoSphere

**AGRICULTURE**  
IBM Watson Analytics, QRM, IBM InfoSphere

**INDUSTRIAL**  
IBM Watson Analytics, QRM, IBM InfoSphere

**OTHER**  
IBM Watson Analytics, QRM, IBM InfoSphere

### OPEN SOURCE

**FRAMEWORKS**  
TensorFlow, PyTorch, IBM Watson Analytics, QRM, IBM InfoSphere

**FORMAT**  
JSON, XML, IBM Watson Analytics, QRM, IBM InfoSphere

**QUERY / DATA FLOW**  
Apache Airflow, IBM Watson Analytics, QRM, IBM InfoSphere

**DATA ACCESS**  
IBM Watson Analytics, QRM, IBM InfoSphere

**DATABASES**  
MySQL, PostgreSQL, IBM Watson Analytics, QRM, IBM InfoSphere

**ORCHESTRATION**  
IBM Watson Analytics, QRM, IBM InfoSphere

**INFRA-STRUCTURE**  
IBM Watson Analytics, QRM, IBM InfoSphere

**DATA OPS**  
IBM Watson Analytics, QRM, IBM InfoSphere

**STREAMING & MESSAGING**  
IBM Watson Analytics, QRM, IBM InfoSphere

**STAT TOOLS & LANGUAGES**  
IBM Watson Analytics, QRM, IBM InfoSphere

**ML OPS & INFRA**  
IBM Watson Analytics, QRM, IBM InfoSphere

**AI / MACHINE LEARNING / DEEP LEARNING**  
IBM Watson Analytics, QRM, IBM InfoSphere

**DATA MARKETPLACES & DISCOVERY**  
IBM Watson Analytics, QRM, IBM InfoSphere

**FINANCIAL & ECONOMIC DATA**  
IBM Watson Analytics, QRM, IBM InfoSphere

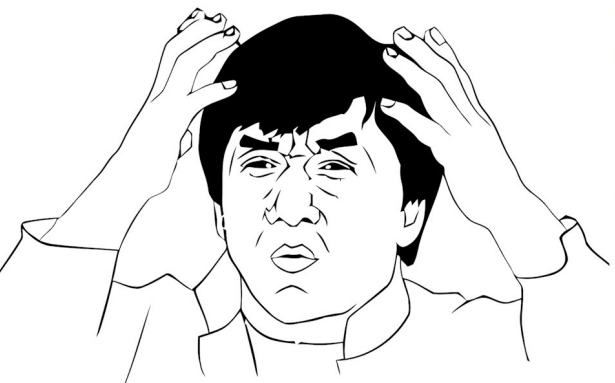
**AIR / SPACE / SEA**  
IBM Watson Analytics, QRM, IBM InfoSphere

**PEOPLE / ENTITIES**  
IBM Watson Analytics, QRM, IBM InfoSphere

**LOCATION INTELLIGENCE**  
IBM Watson Analytics, QRM, IBM InfoSphere

**OTHER**  
IBM Watson Analytics, QRM, IBM InfoSphere

**DATA SERVICES**  
IBM Watson Analytics, QRM, IBM InfoSphere



# Распространенные MLOps-инструменты



## OPERATIONALIZATION

### MODEL MONITORING

arize EVIDENTLY AI fiddler LOSSWISE unravel

### MODEL DEPLOYMENT/SERVING

BENTOML Kubeflow SELDON TensorFlow  
TensorFlow Serving

## MODELING

### FEATURE ENGINEERING

dotData FEAST Featuretools RASGO tsfresh

### MODEL VERSIONING

DVC mlflow ModelDB neptune.ai

### EXPERIMENT TRACKING

comet neptune.ai Snorkel TensorFlow

### HYPERPARAMETER OPTIMIZATION

HYPEROPT SCIKIT-OPTIMIZE SIGOPT

## DATA MANAGEMENT

### DATA LABELING

doccano iMerit Labelbox Prodigy

### DATA STORAGE & VERSIONING

comet DVC dolt lakeFS Pachyderm QRI

## END-TO-END MLOPS

Amazon SageMaker Azure Machine Learning CLEAR ML CLUSTER CLUSTER databricks DataRobot DOMINO H2O.ai iguazio Weights & Biases Valohai Vertex AI

# Основные инструменты



15

**Версионирование и отслеживание экспериментов:** MLflow, DVC (Data Version Control), WANDB

**CI/CD пайплайны:** Jenkins, GitLab CI, CircleCI, Kubeflow Pipelines

**Контейнеризация и оркестрация:** Docker, Kubernetes

**Мониторинг:** Prometheus, Grafana, AWS CloudWatch, Azure Monitor

**DataOps:** Apache Airflow, Prefect, Kafka, DBT

**ML Платформы:** Weights & Biases, Neptune.ai, ClearML.

**Сервинг моделей:** TensorFlow Serving, TorchServe, Flask, FastAPI, KServe, Nvidia Triton, Nuclio (Serverless)

# Различие в инфраструктуре при обучении модели и для инференса

	Обучение	Инференс
Цель	Улучшение качества модели	Меньшая задержка на ответ
Вычислительные ресурсы	Необходимы большие ресурсы, в т.ч. GPU	Ресурсов требуется меньше, GPU тоже меньше объема
Загрузка	Батч загрузка, много длинных задач	Запрос/Ответ, может быть батч. Оптимизировано под быстрый ответ
Необходимость масштабирования	Большие кластеры для обучения крупных моделей	Автомасштабирование при увеличении нагрузки
Задержка	-	Адекватное для человека время на ответ
Дисковое пространство	Большие данные, озера данных	Минимальные объемы (модель, конфиг, логи, дополнительный код)
Отказоустойчивость	-	Высокая избыточность, важность безотказной работы
Мониторинг	Mlflow и подобные для трекинга экспериментов	Постоянно следить за качеством модели, оценивать data drift и performance degradation



# Деплоймент моделей (инференс)

## По типу выполнения:

### Онлайн инференс:

- Специализированные облачные решения (AWS SageMaker, GCP Cloud Run, Selectel Inference Platform и др.)
- Или развернутые там же opensource инструменты

ОБЛАКО

- Существует множество решений, подходящих под разные технологии, но нет единого стандартного
- Необходимо подбирать инструмент под конкретную технологию
- С развитием LLM начали появляться специализированные решения для инференса с возможностью запуска (квантизованных) и обычных моделей ollama, vllm и т.д.

СВОЯ ИНФРАСТРУКТУРА

## По типу:

- REST/GRPC
- Model as a Service (MAS)
- Model-on-Demand
- ...

### Оффлайн инференс:

- Батч обработка
- REST/GRPC

# Деплоймент Model as a Service (REST/GRPC)



```
from fastapi import FastAPI
from pydantic import BaseModel

class ModelRequest(BaseModel):
    name: str
    price: float

app = FastAPI()

@app.post("/predict/")
async def next_best_offer(request: ModelRequest):
    return modelService.getModel().predict(request)
```

+

Ds-way подойдет для быстрого прототипа решения

-

Нет поддержки cloud native и интеграции с k8s, как следствие, решение не масштабируется

-

Непонятно, как будет работать под нагрузкой?

-

Быстрее будет собирать запросы в батчи или отправлять в модель только по одному?

-

Если потребуется добавить новый функционал (например, авторизацию и аутентификацию), нужно внести изменения в код. Изменения могут быть довольно частыми

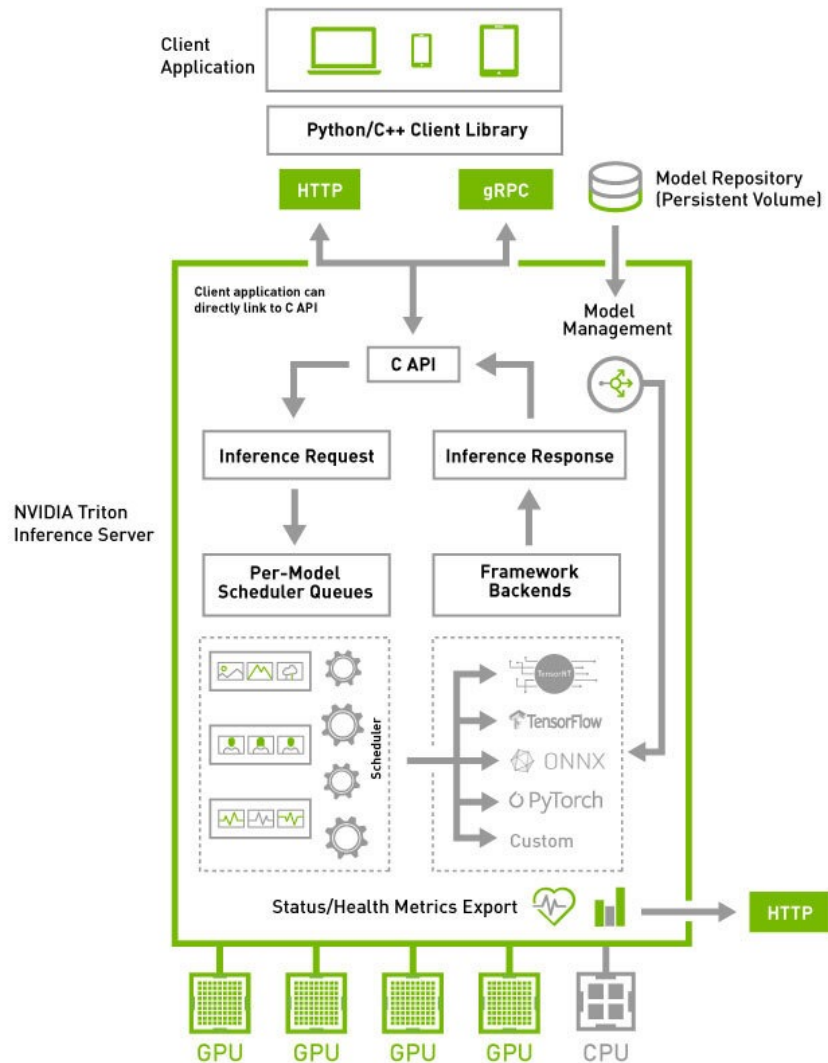
-

Нет мониторинга качества и непонятно, когда модель станет деградировать

-

Может меняться как код, так и сама модель, придется все переписывать заново

# Inference Servers



## Triton Inference Server

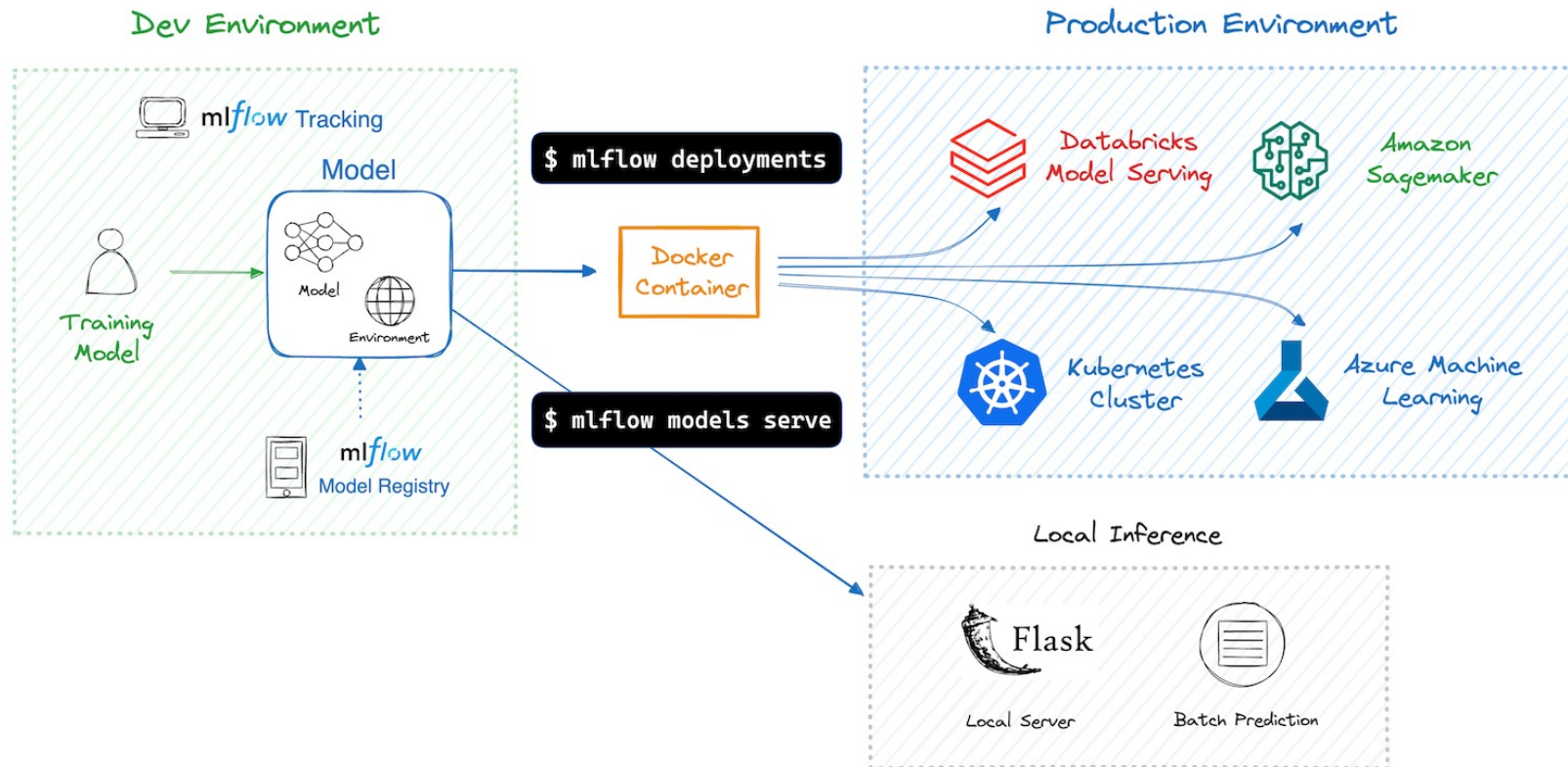
- + Одновременное выполнение моделей
- + GPU/CPU
- + Адаптивный (динамический) батчинг
- + Поддержка практически всех известных бэкендов (TensorRT, TensorFlow, PyTorch, ONNX, OpenVINO, Python и другие)
- + Горячая замена моделей
- + Много дополнительных функций
- Весьма сложен

Другие:

- TFServing - tensorflow
- TorchServe - pytorch

Fall into ML '24

# MLFlow Serving



## Локальный инференс

- Для тестирования
- Не подходит под прод нагрузку
- + Быстро и просто развернуть

## Прод инференс

- + Лишен недостатков локального
- + Асинхронный, отдельный пул обработчиков, можно одновременно сервить несколько моделей

Fall into ML '24

# Model Serving. Kserve



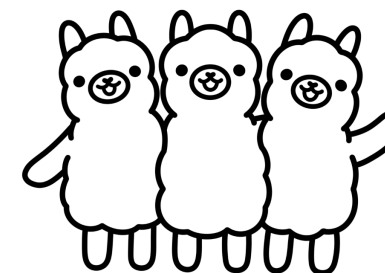
- Поддерживает большинство известных райнтаймов TF Serving, TorchServe, Triton Inference Server
- Cloud native
  - Масштабирование на CPU/GPU
  - Управление версиями
  - Батчинг
  - Логирование (запрос/ответ)
  - Управление трафиком
  - Метрики

и др.

# Специализированные решения

LLM - vLLM, Ollama, llama.cpp

- llama3.2 1b, 3b
- llama3.1 8b,70b,450b
- gemma2 2b,9b,27b
- qwen2.5 0.5...72b
- mistral-nemo 12b
- mistral
- mixtral
- llava 7b,13b,34b



<https://ollama.com/library>

Fall into ML  
'24

# ЗАГОРУЛЬКИН ДМИТРИЙ ЭДУАРДОВИЧ

Заместитель директора Центра  
стратегической аналитики и больших данных  
ИСИЭЗ НИУ ВШЭ

Заместитель руководителя системы  
интеллектуального анализа больших данных  
iFORA

[dzagorulkin@hse.ru](mailto:dzagorulkin@hse.ru)



Сайт iFORA



iFORA в Telegram



iFORA-экспрессы

Хотите у нас работать  
или пройти  
стажировку?

**Сканируйте QR-код**

