

@AxSaucedo

PyCon HK 2020

Real Time Machine Learning with Python

Alejandro Saucedo | as@seldon.io

Twitter: [@AxSaucedo](https://twitter.com/AxSaucedo)

Hello, my name is Alejandro

@AxSaucedo



Alejandro Saucedo

Engineering Director
Seldon Technologies

Chief Scientist
The Institute for Ethical AI & ML

Head of Solutions Eng & Sci
Eigen Technologies

Software Engineer & DevX Lead
Bloomberg LP

Seldon: OSS Production ML Deployment

@AxSaucedo

1. Package

Create REST or gRPC dockerized microservice .

2. Describe Deployment

Create/update Kubernetes resource manifest for deployment graph.

3. Deploy

Manage and analyze the performance of live deployments.

2. Seldon Deploy

(UI, Collaboration, Control, Audit)

MAB
(Multi-Arm Bandits)

Outlier Detection

Explanation

Bias Detection

1. Seldon Core

(runtime ML graph engine)

Microservices - Istio service mesh (optional)



kubernetes

The Institute for Ethical AI & Machine Learning

@AxSaucedo



The Institute for Ethical AI & Machine Learning

We are a UK-based think tank that brings together technologists, policymakers & academics to develop standards & frameworks for Data Governance & Machine Learning.

We are part of the LFAI

@AxSaucedo



Trusted AI Committee



Today

@AxSaucedo

- **Conceptual intro to stream processing**
- **Machine learning for real time**
- **Tradeoffs across tools**
- **Hands on use-case**

Real Time Reddit Processing

@AxSaucedo

- Real time ML model for reddit comments
- 200k comments for training model
- /r/science comments removed by mods

**We will be fixing the front page
of the internet**

A trip to the ~~past~~ present: ETL

@AxSaucedo



E - Extract
T - Transform
L - Load

Variations

@AxSaucedo

- **ETL - Extract Transform Load**
- **ELT - Extract Load Transform**
- **EL - Extract Load**
- **LT - Load Transform**
- **WTF - LOL**

Specialised Tools

@AxSaucedo



EL

ETL

ELT

Nifi

Oozie

Elasticsearch

Flume

Airflow

Data Warehouse

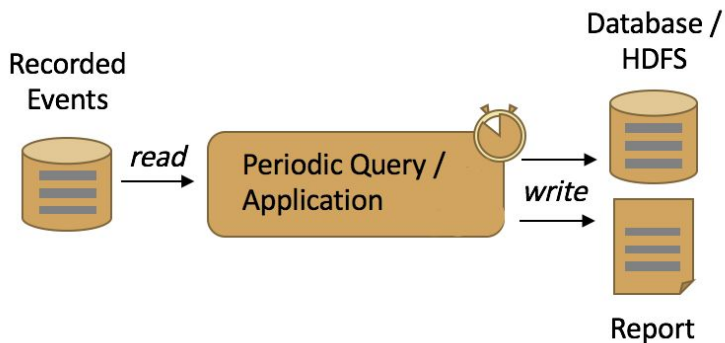
...

Jupyter notebook?

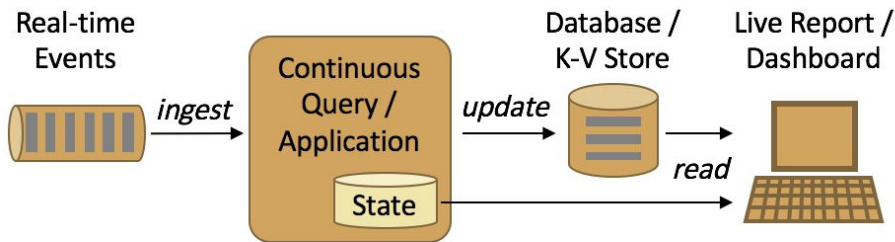
Batch VS Streaming

@AxSaucedo

Batch analytics



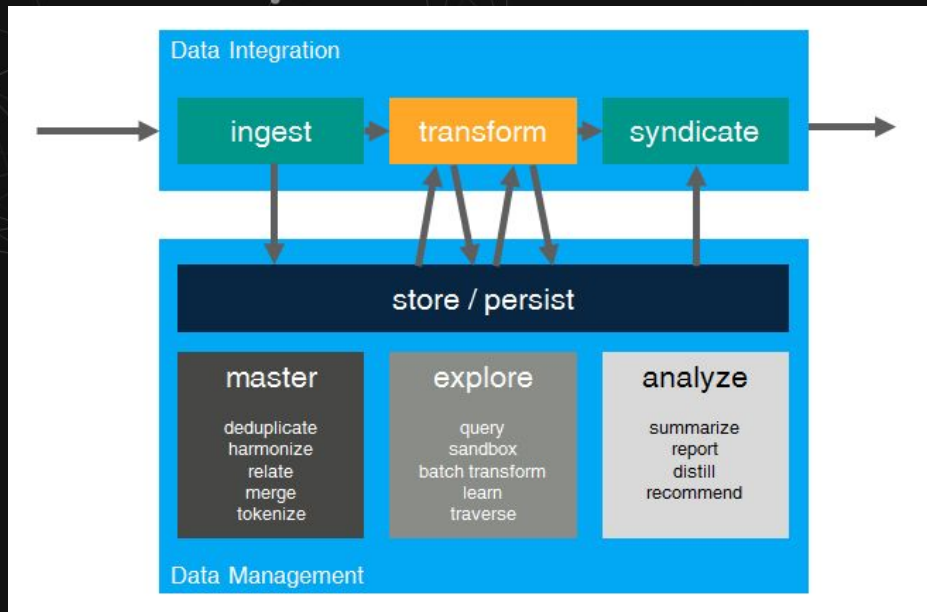
Streaming analytics



The spectrum of data processing

Batch ~~VS~~ AND Streaming

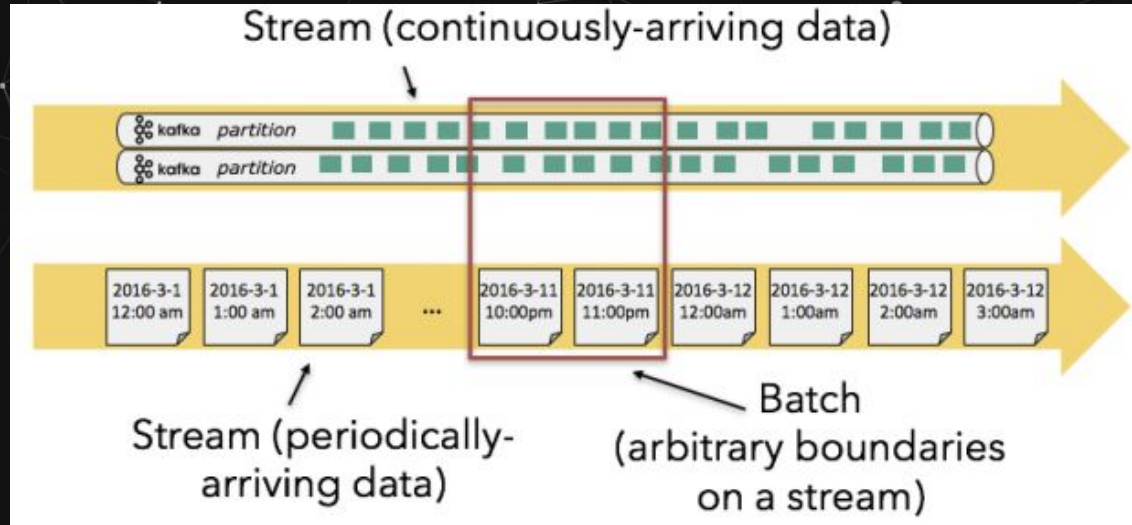
@AxSaucedo



The right tool for the challenge

Unifying Worlds

@AxSaucedo



Massive drive on converging worlds

Streaming Concepts: Windows

@AxSaucedo

- Tumbling windows



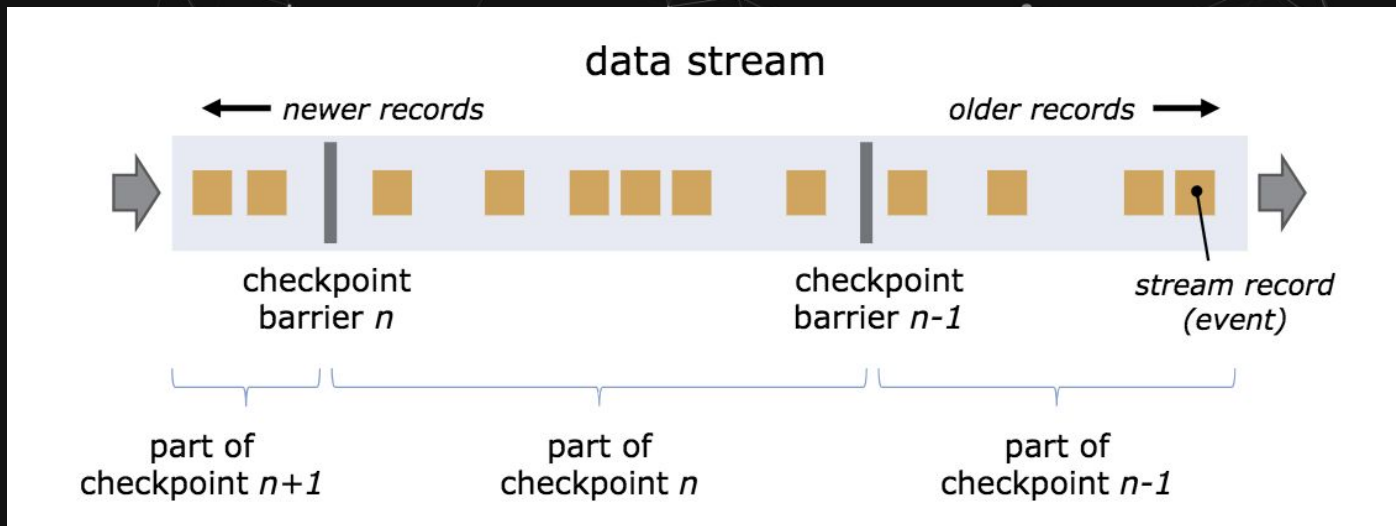
- Sliding windows



Processing of batches in real time

Streaming Concepts: Checkpoint

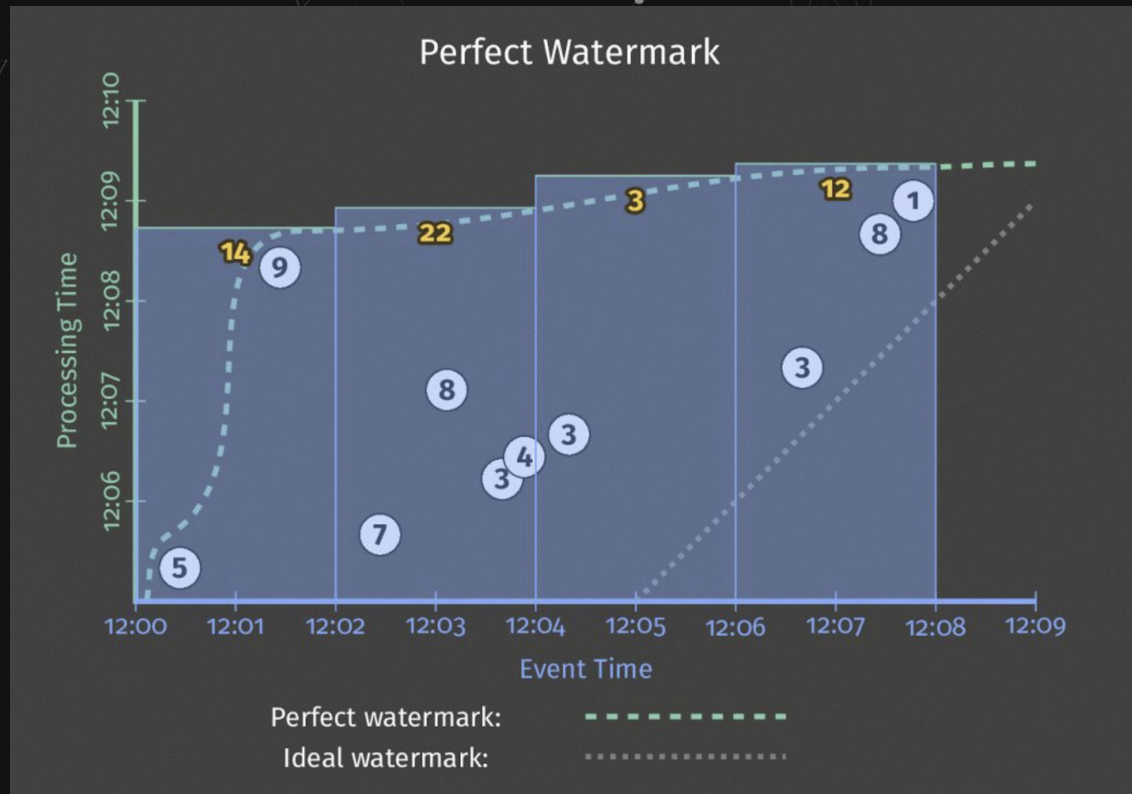
@AxSaucedo



Keeping track of stream progress

Streaming Concepts: Watermark

@AxSaucedo



Considering
data that
comes late in
windows and
stream
batches

Some Stream Processing Tools

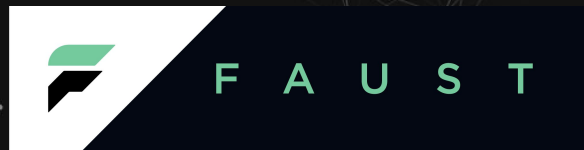
@AxSaucedo

- **Flink (Multiple Languages)**
- **Kafka Streams (Multiple Languages)**
- **Spark Stream (Multiple Languages)**
- **Faust (Python)**
- **Apache Beam (Python)**

Today we're using

@AxSaucedo

Stream Processing



ML Serving



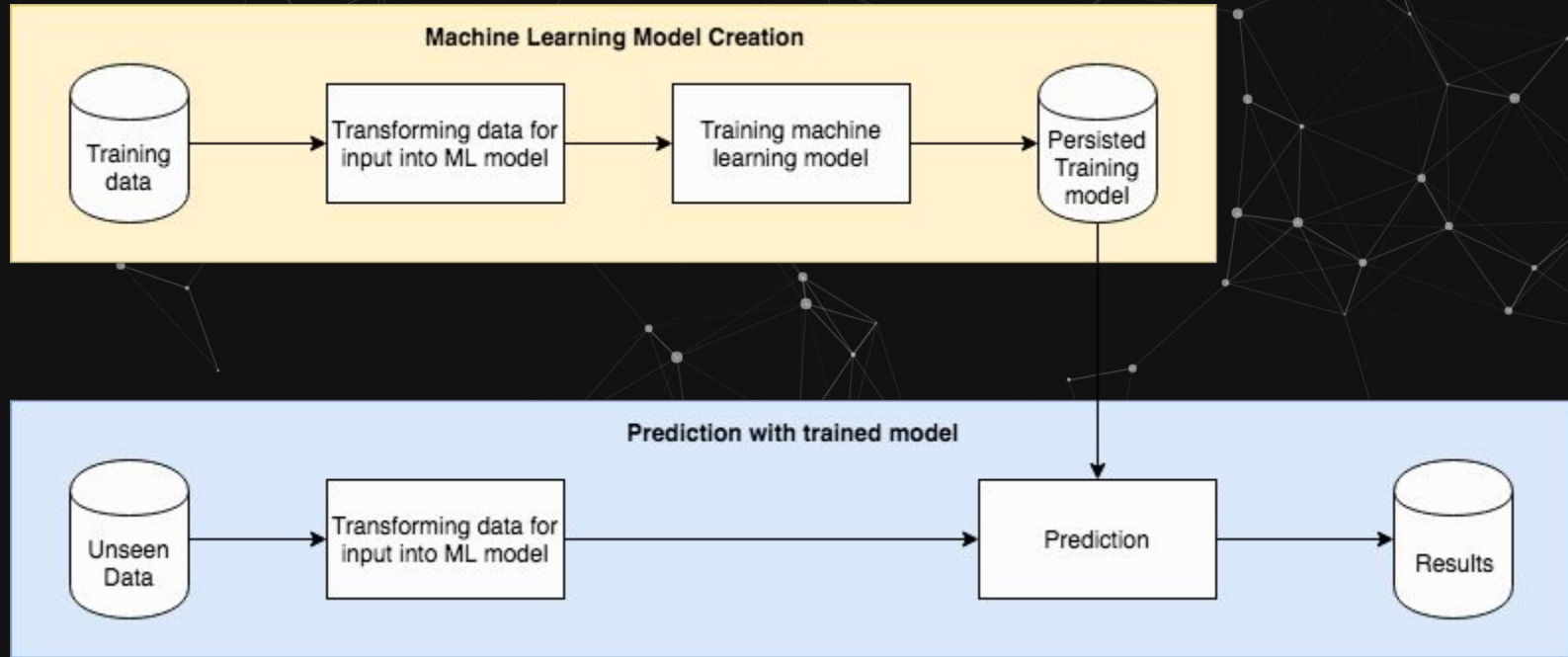
ML Training

spaCy



Machine Learning Workflow

@AxSaucedo



Model Training

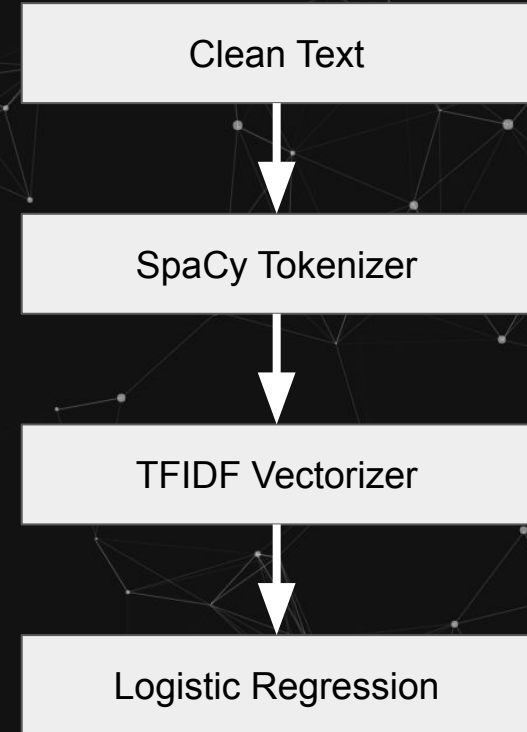
@AxSaucedo

```
clean_text_transformer = CleanTextTransformer()

spacy_tokenizer = SpacyTokenTransformer()

tfidf_vectorizer = TfidfVectorizer(
    min_df=3,
    max_features=1000,
    preprocessor=lambda x: x, tokenizer=lambda x: x,
    token_pattern=None,
    ngram_range=(1, 3), use_idf=1, smooth_idf=1,
    sublinear_tf=1)

lr_model = LogisticRegression(C=1.0, verbose=True)
```



"You are a DUMMY!!!!!"

Model Training

@AxSaucedo

"You are dummy"

[PRON, IS, DUMB]

[1000, 0100, 0010]

[1]

```
x_train_clean = \
    clean_text_transformer.transform(x_train)

x_train_tokenized = \
    spacy_tokenizer.transform(x_train_clean)

tfidf_vectorizer.fit(
    x_train_tokenized[TOKEN_COLUMN].values)

x_train_tfidf = \
    tfidf_vectorizer.transform(
        x_train_tokenized[TOKEN_COLUMN].values)

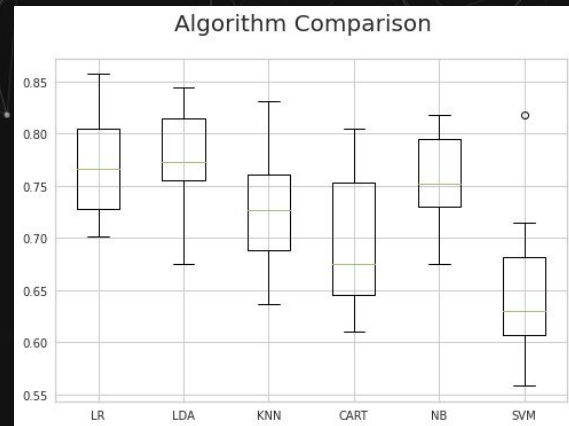
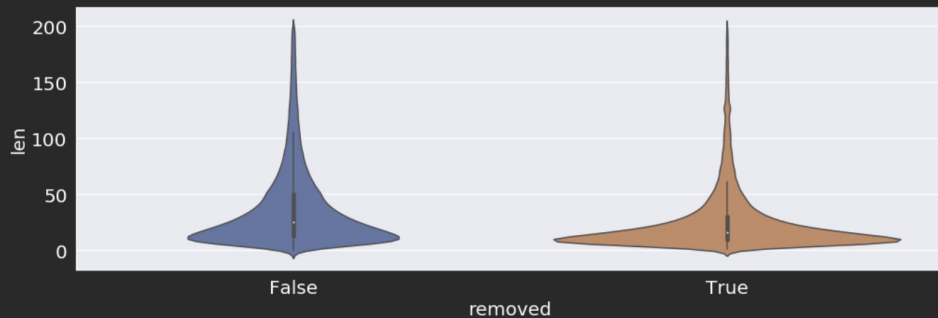
lr_model.fit(x_train_tfidf, y_train)

pred = lr_model.predict(x_test_tfidf)
```

More on EDA & Model Evaluatio

@AxSaucedo

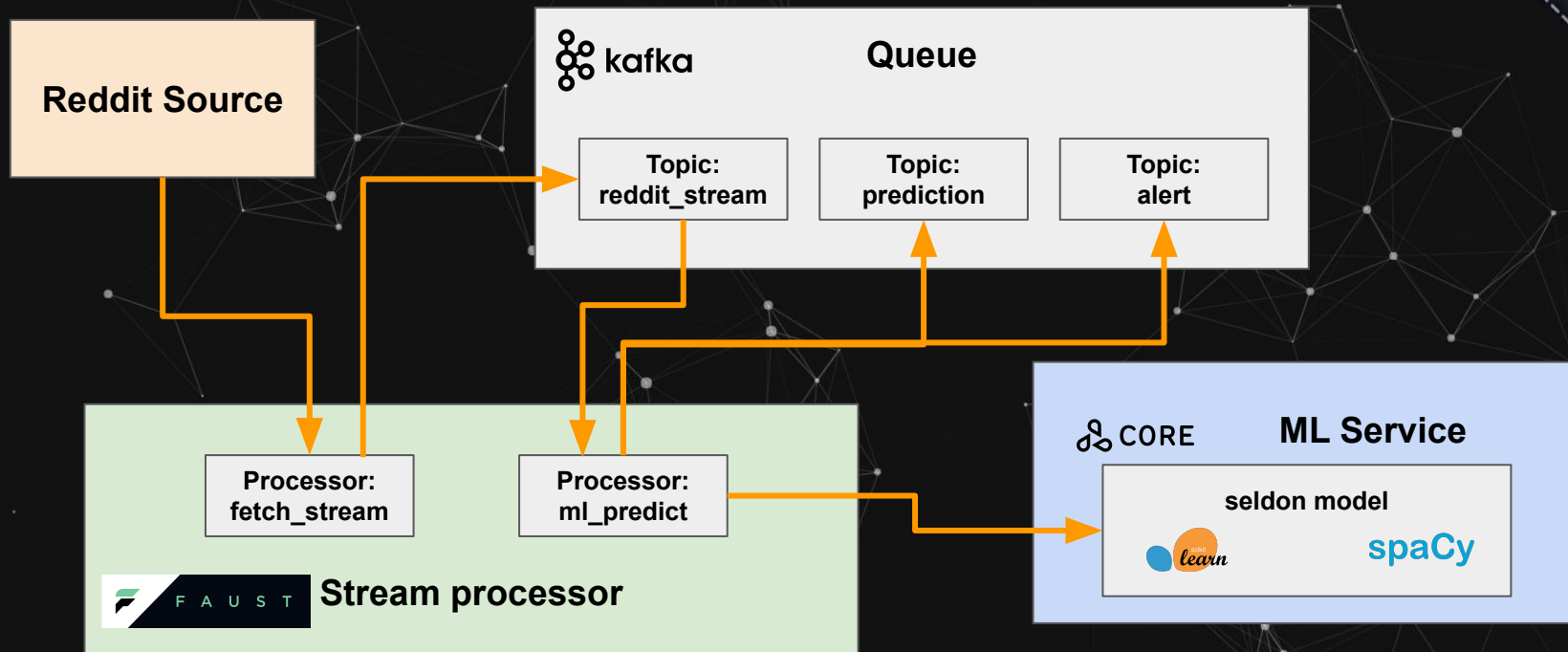
```
plt.show()
sns.catplot(x="removed", y="len", data=df_extended[df_len<200], kind="violin", aspect=20/7)
```



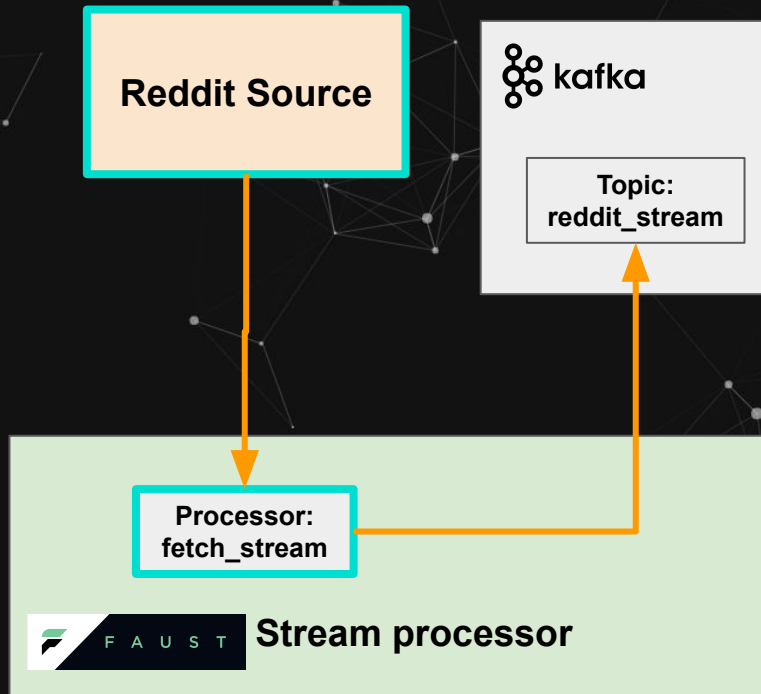
<https://github.com/axsaucedo/reddit-classification-exploration/>

Overview of Components

@AxSaucedo



Generating comments



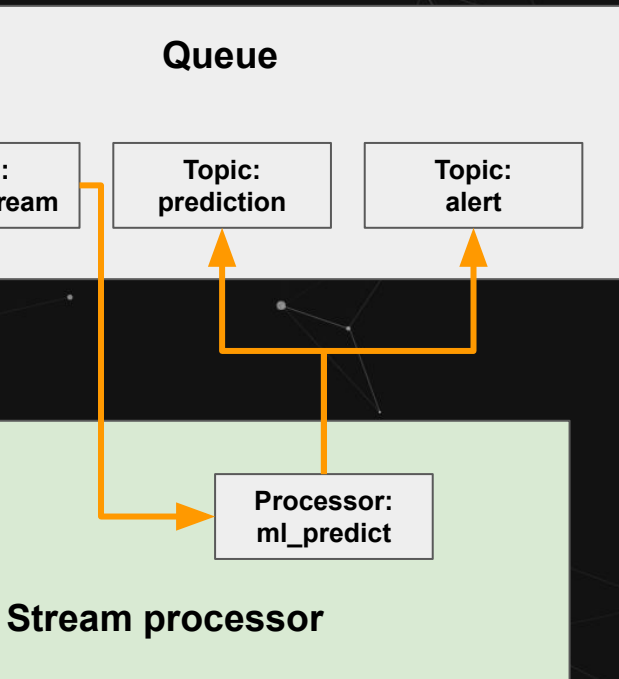
```
@app.timer(0.1)
async def generate_reddit_comments():
    reddit_sample = await fetch_reddit_comment()

    reddit_data = {
        "id": reddit_sample["id"].values[0],
        "score": int(reddit_sample["score"].values[0]),
        ... # Cut down for simplicity
    }

    await app.topic("reddit_stream").send(
        key=reddit_data["id"],
        value=reddit_data)
```

ML Stream Processing Step

@AxSaucedo



```
@app.agent(app.topic("reddit_stream"))
async def predict_reddit_content(tokenized_stream):
    async for key, comment_extended in tokenized_stream.items():
        tokens = comment_extended["body_tokens"]

        probability = seldon_prediction_req(tokens)

        data = {
            "probability": probability,
            "original": comment_extended["body"]
        }

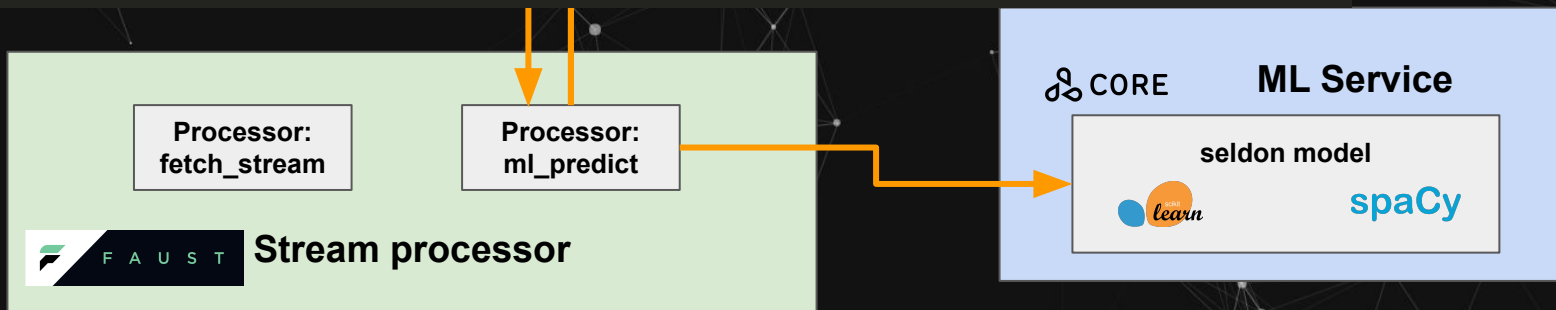
        await app.topic("reddit_prediction").send(
            key=key,
            value=data)

        if probability > MODERATION_THRESHOLD:
            await reddit_mod_alert_topic.send(
                key=key,
                value=data)
```


ML Model Request Step

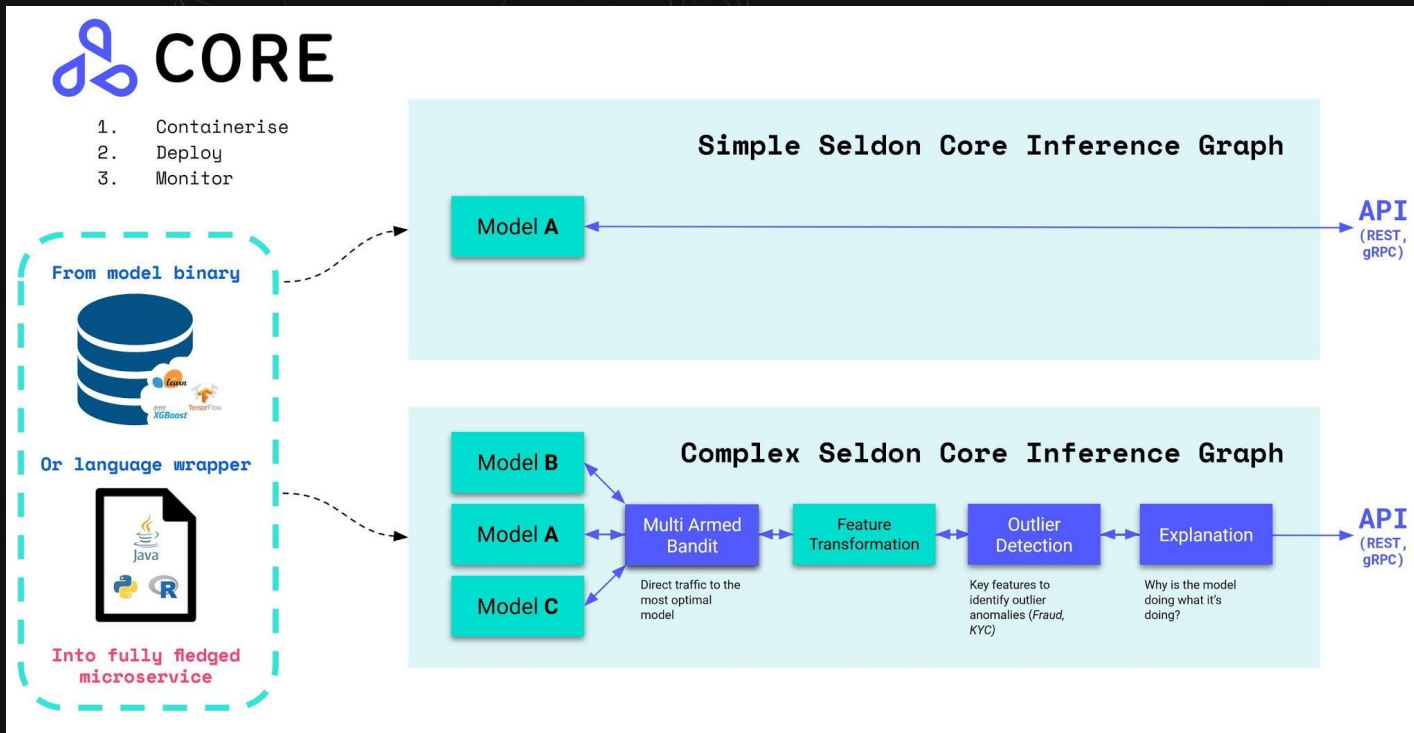
@AxSaucedo

```
sc = SeldonClient(  
    gateway_endpoint="istio-ingress.istio-system.svc.cluster.local",  
    deployment_name="reddit-model",  
    namespace="default")  
  
def seldon_prediction_req(tokens):  
    data = np.array(tokens)  
    output = sc.predict(data=data)  
    return output.response["data"]["ndarray"]
```



Overview of Seldon Model Serving

@AxSaucedo



Wrapping ML models for Serving with Seldon

```
import dill

from ml_utils import CleanTextTransformer, SpacyTokenTransformer

class RedditClassifier:

    def __init__(self):
        self._clean_text_transformer = CleanTextTransformer()
        self._spacy_tokenizer = SpacyTokenTransformer()

        with open('tfidf_vectorizer.model', 'rb') as model_file:
            self._tfidf_vectorizer = dill.load(model_file)

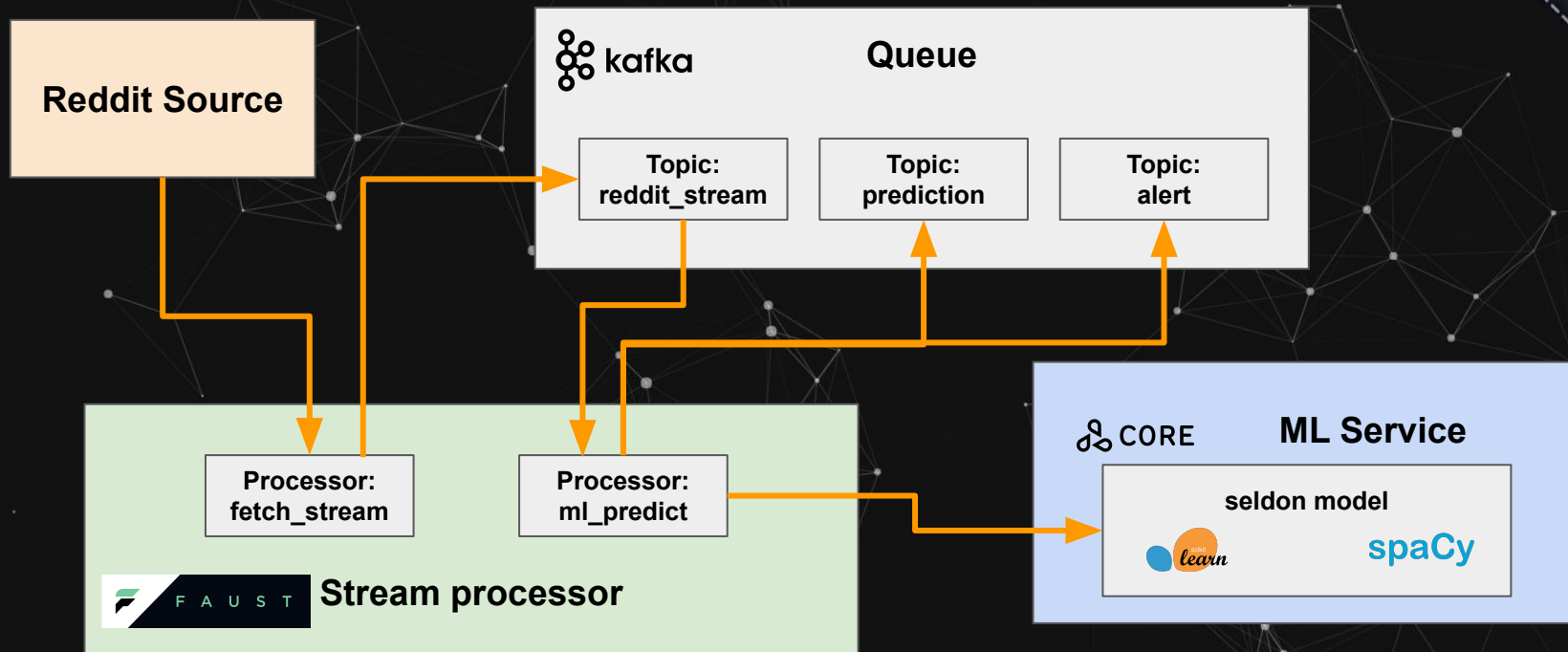
        with open('lr.model', 'rb') as model_file:
            self._lr_model = dill.load(model_file)

    def predict(self, X, feature_names):
        clean_text = self._clean_text_transformer.transform(X)
        spacy_tokens = self._spacy_tokenizer.transform(clean_text)
        tfidf_features = self._tfidf_vectorizer.transform(spacy_tokens)
        predictions = self._lr_model.predict_proba(tfidf_features)
        return predictions
```

@AxSaucedo

Overview of Components

@AxSaucedo



Recap of Today

@AxSaucedo

- **Conceptual intro to stream processing**
- **Machine learning for real time**
- **Tradeoffs across tools**
- **Hands on use-case**

Seldon Technologies

Real Time Machine Learning with Python

Alejandro Saucedo | as@seldon.io

[@AxSaucedo](https://twitter.com/AxSaucedo)