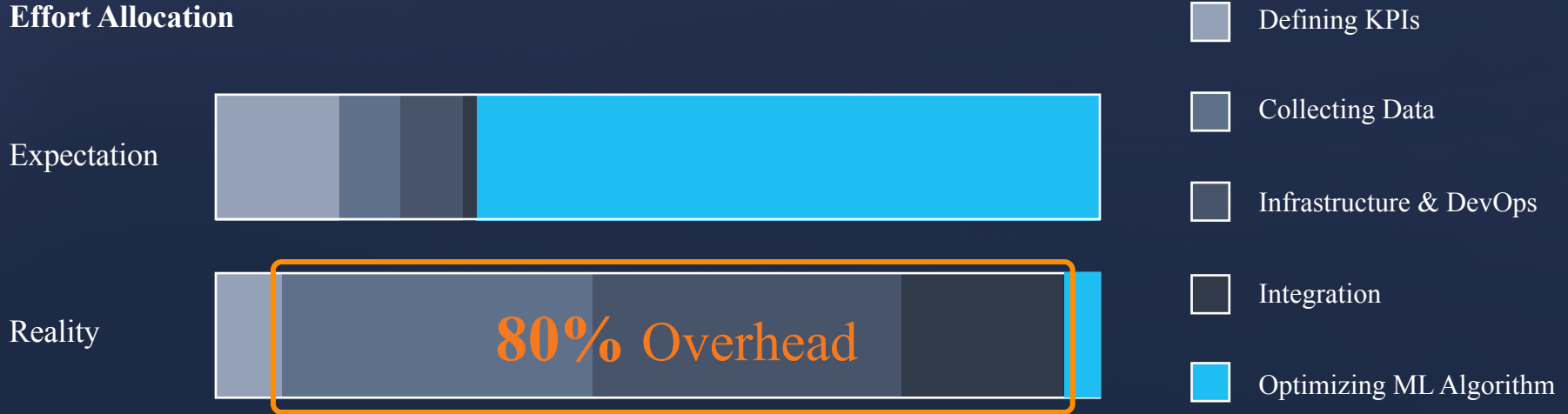# Nuclio : KubeFlow Serverless's  component

Orit Nissan-Messing, VP R&D, Iguazio

# Data Science Teams Don't Do Data Science

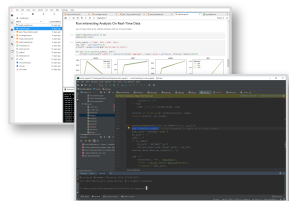**Effort Allocation**

Expectation

Reality

**80% Overhead**

Defining KPIs

Collecting Data

Infrastructure & DevOps

Integration

Optimizing ML Algorithm

**Source: Google Developers Launchpad**

**The need: Simpler Solutions, Better Data Integration**

1

iguazio

**Develop/Experiment**   **Package**   **Scale-out**   **Tune**   **Instrument**   **Automate**

- Dependencies
- Parameters
- Run scripts
- Build

- Load-balance
- Data partitions
- Model distribution
- Hyper params

- Parallelism
- GPU support
- Query tuning
- Caching

- Monitoring
- Logging
- Versioning
- Security

- CI/CD
- Workflows
- Rolling upgrades
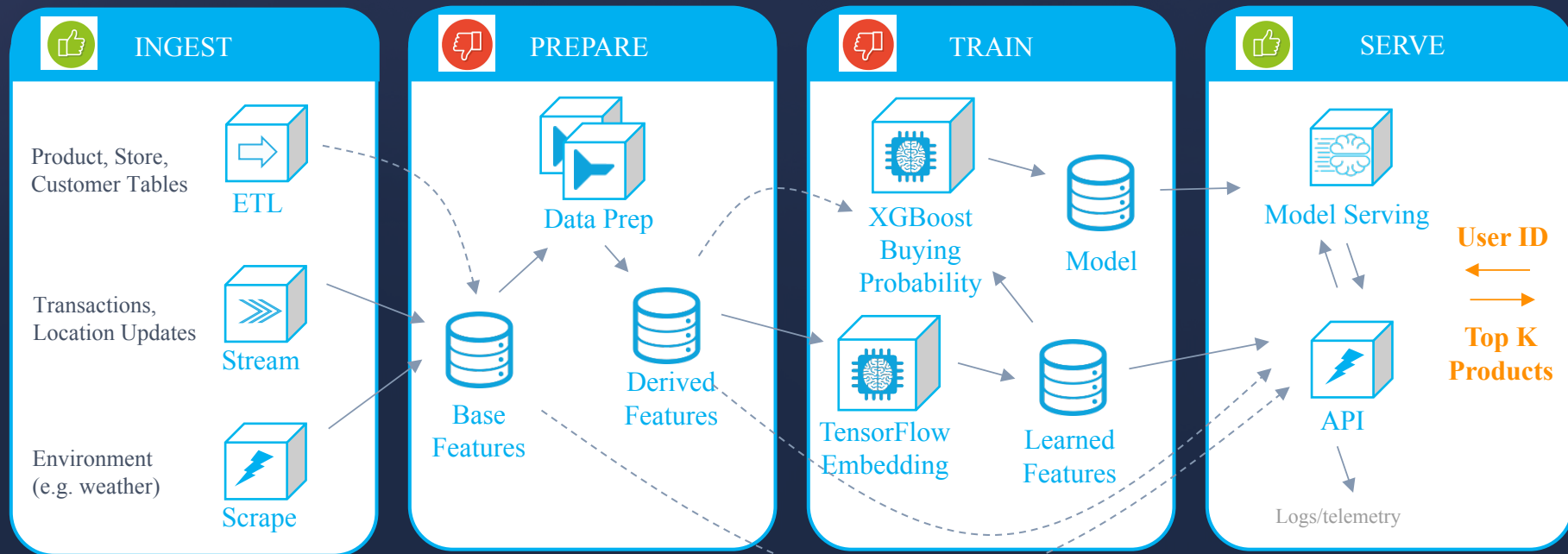- A/B testing

**Weeks** with one data scientist

**Months** with a large team of developers, scientists, data engineers and DevOps

Automate DevOps to Deploy Projects in

One Week as Opposed to Months!

2

iguazio

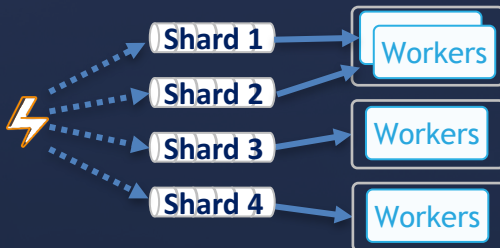Example: Real-time Product Recommendations

# Nuclio: Taking Serverless to Data Intensive Apps
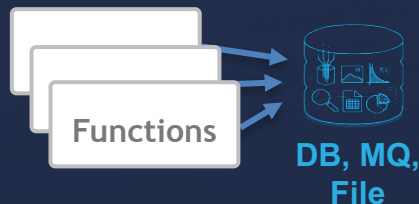
## Extreme Performance



- Non-blocking, parallel
- Zero copy, buffer reuse
- Up to 400K events/sec/proc
- **GPU** optimizations

## Advanced Data & AI Features



- Auto-rebalance, checkpoints
- Any source: Kafka, NATS, Kinesis, event-hub, iguazio, pub/sub, RabbitMQ, Cron, ..
- NVIDIA Rapids integration

## Statefulness



- Data bindings
- Shared volumes
- Context cache

## Natively integrated with Kubeflow and Jupyter Notebooks

iguazio

# Using Nuclio to Accelerate ETL and Streaming

**Simple code!  Automated DevOps !  Any Source!**
(e.g. read JSON Stream + aggregate + dump to Parquet)

```python
def init_context(context):
    os.makedirs(sink, exist_ok=True)

def handler(context, event):
    add_log_to_batch(context, event.body)

    if len(batch) > batch_len:
        df = _batch_to_df(context)
        if not df.empty:
            df = df.groupby(['log_ip']).agg({'feconn':'mean',
                                             'beconn':'mean',
                                             'time_backend_response':'max',
                                             'time_backend_response':'mean',
                                             'time_queue':'mean',
                                             'time_duration': 'mean',
                                             'time_request': 'mean',
                                             'time_backend_connect':'mean'
                                            })
        df_to_parquet(df)
        reset_batch()
```
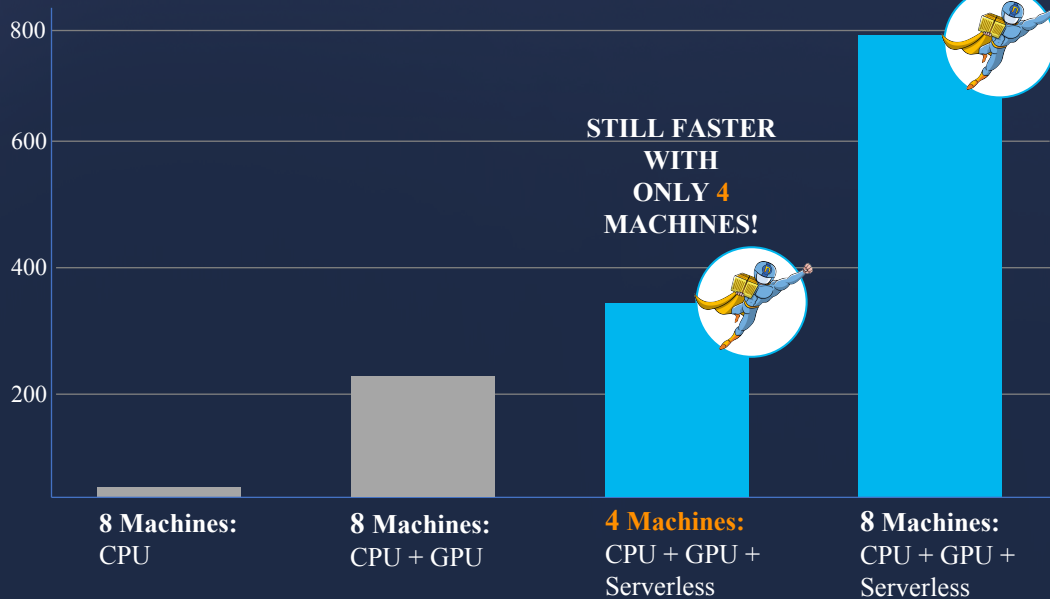
**500 MB/s**

**Simple
Python**

**18 MB/s**

4

iguazio

# Why Not Use Serverless for Training and Data Prep?

*What about Training and data prep ?*

|  | **Serverless Today** | **Data Prep and Training** |
|---|---|---|
| **Task lifespan** | Millisecs to mins | Secs to hours |
| **Scaling** | Load-balancer | Partition, shuffle, reduce, Hyper-params |
| **State** | Stateless | Stateful |
| **Input** | Event | Params, Datasets |

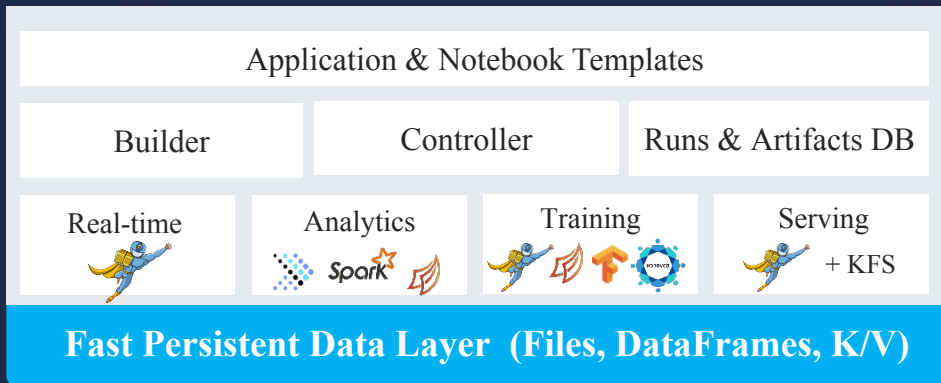**Serverless: resource elasticity (to Zero) and automated deployment and operations**

iguazio

# Introducing Nuclio ML Functions

**Access from your notebook, IDE, or KubeFlow**



**Common APIs & Automation**

**Multiple Engines**

| Application & Notebook Templates |
|---|

| Builder | Controller | Runs & Artifacts DB |
|---|---|---|

| Real-time | Analytics | Training | Serving + KFS |
|---|---|---|---|

**Fast Persistent Data Layer  (Files, DataFrames, K/V)**

**Built-in Artifacts & Runs Tracking**

**Elastic Scaling**

iguazio

# Demo: Fast and Serverless KubeFlow Pipeline



**1** Collect from Multiple Sources

**2** Prepare and Explore

**3** Accelerate and Automate Training

**4** Deploy in One Click

Real-time Data Layer

GPUs

iguazio

**iguazio**

# Thank You

oritn@iguazio.com, www.iguazio.com