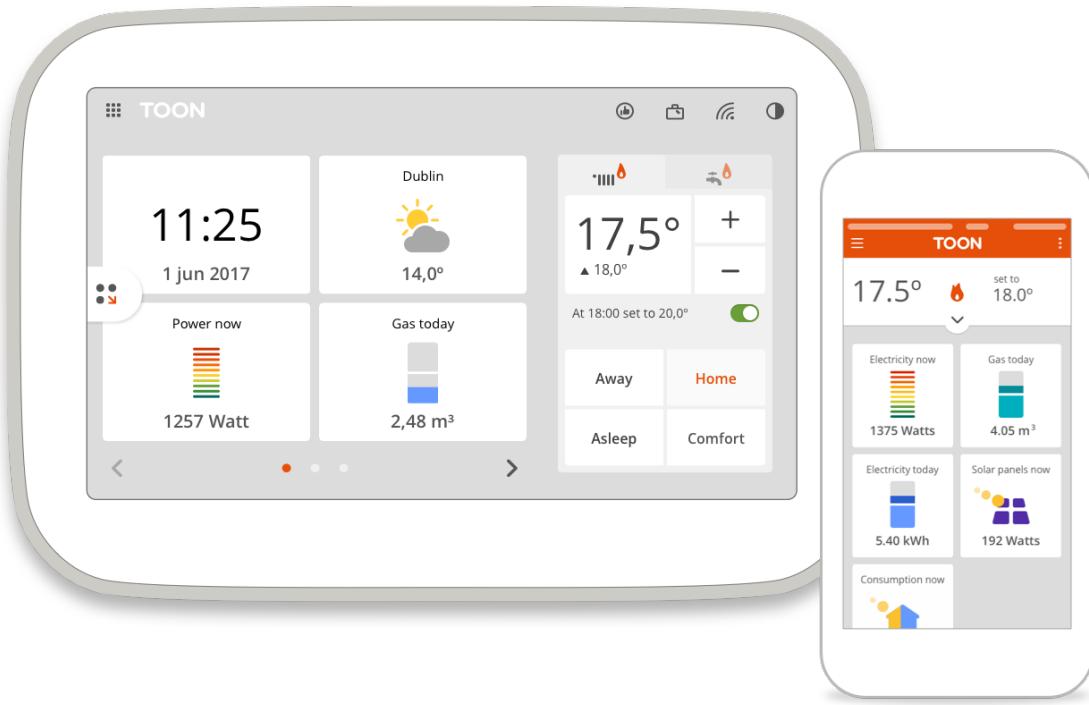


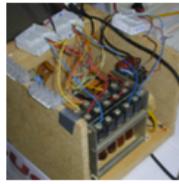
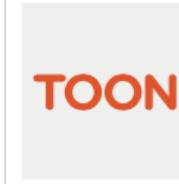


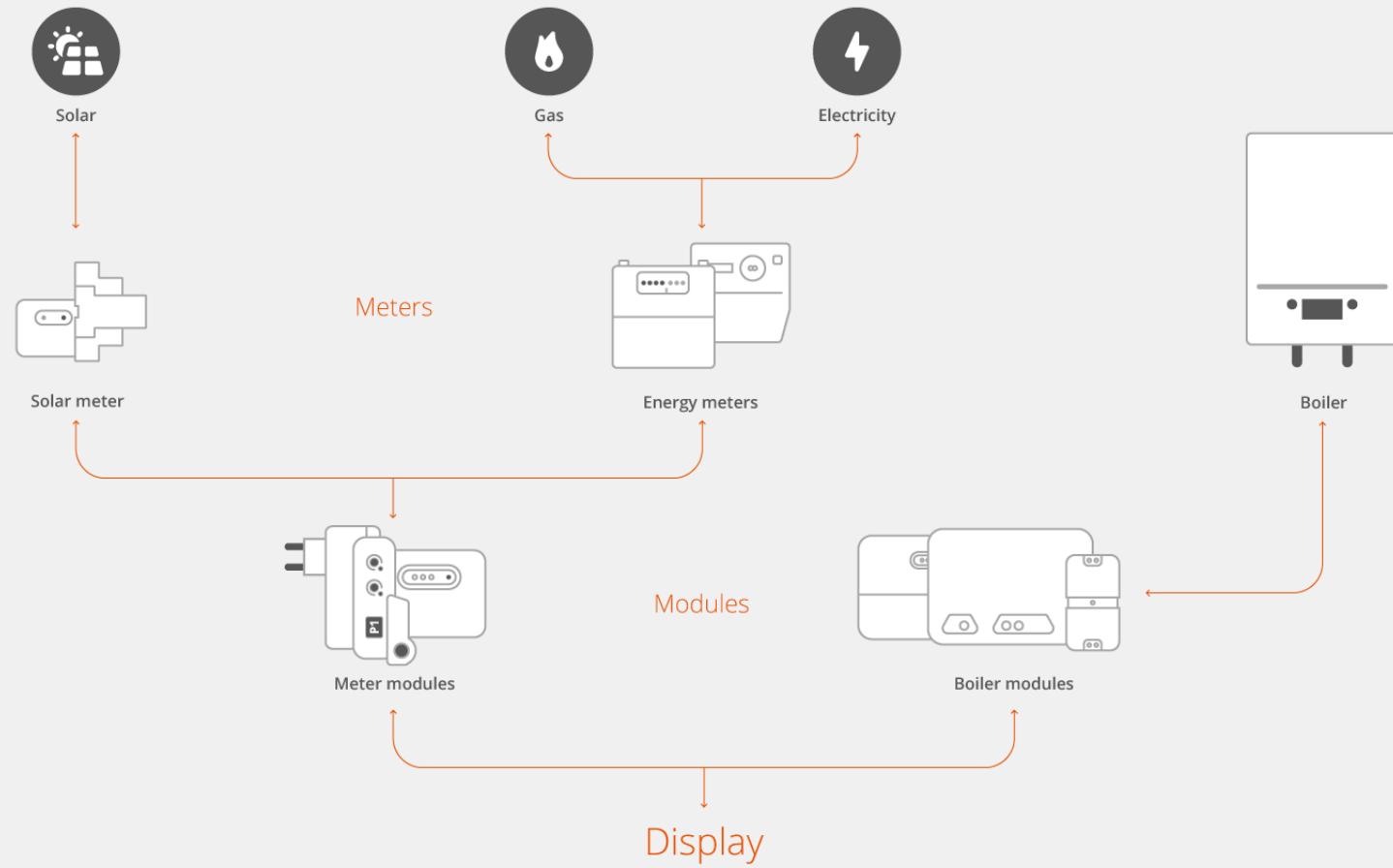
Using Spark in the Cloud: A Devops perspective

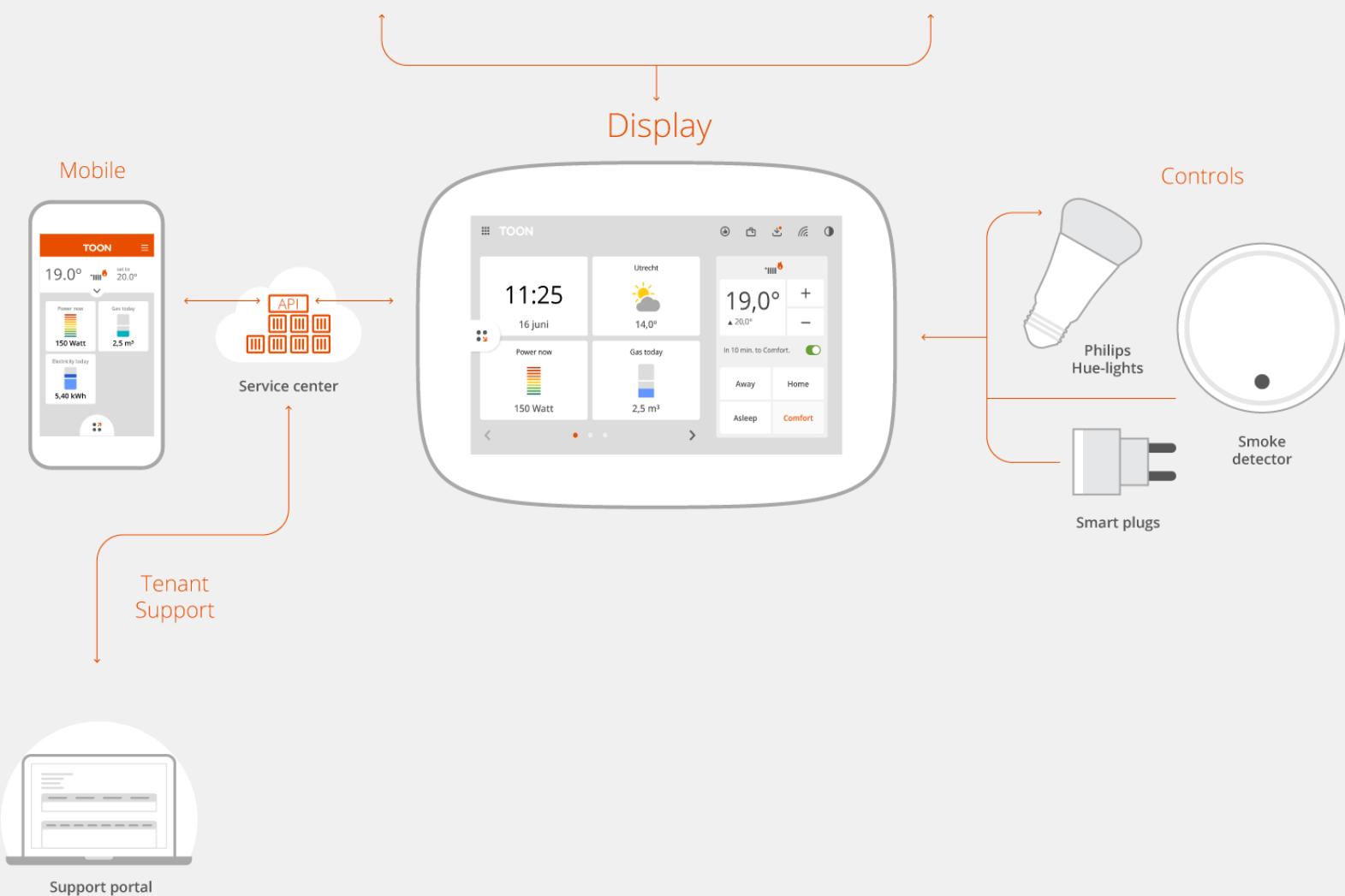
Telmo Oliveira, Toon

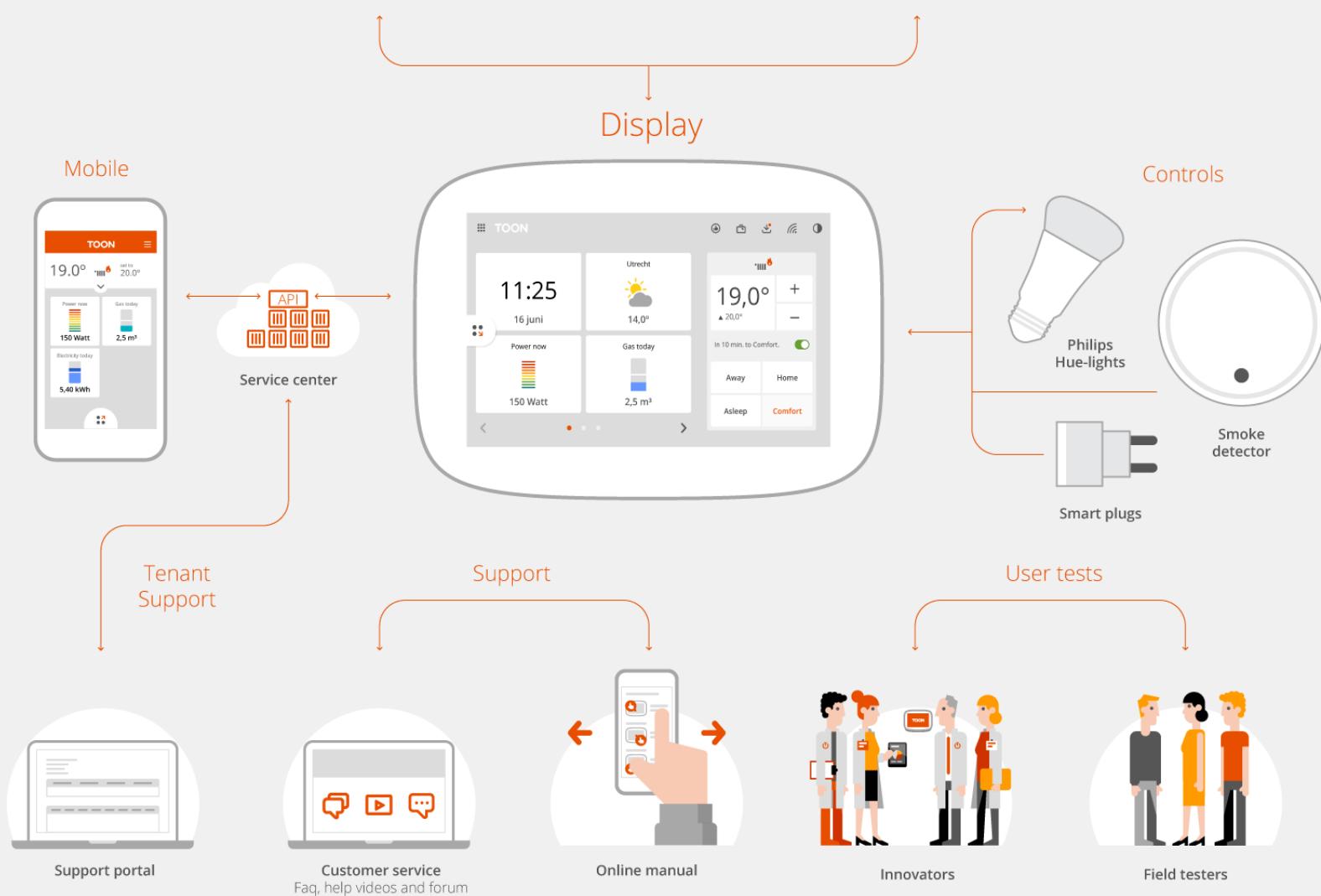


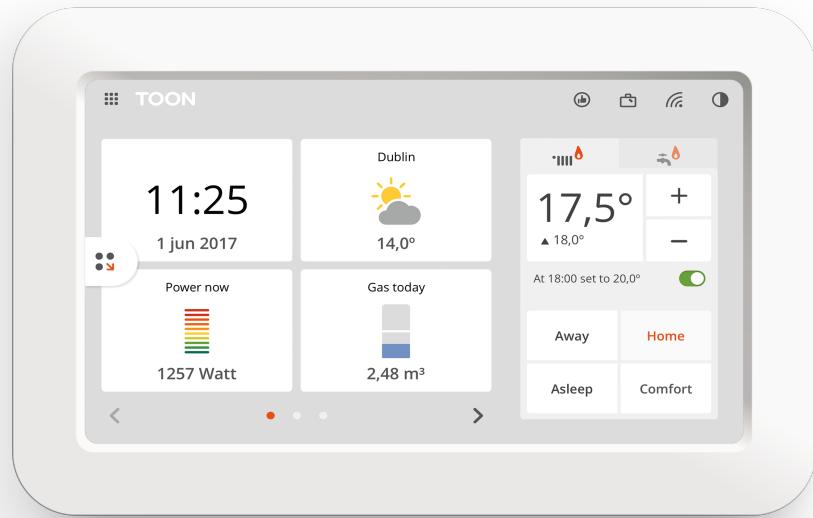
TOON®

2004	Home Automation Europe is founded	Home control alternatives are prototyped	The shift to smart thermostats	Home Automation Europe becomes Quby	Eneco acquires 100% of Quby's shares	Quby's white label product is rebranded internationally as Toon
						
2005						
2006						
2007						
2008						
2009						
2010						
2011						
2012						
2013						
2014						
2015						
2016						
2017						
	The Home Control Box is launched	Home Automation Europe moves into a new office		Home Automation Europe partners with Eneco to create Toon	Quby partners with Engie to create Boxx	

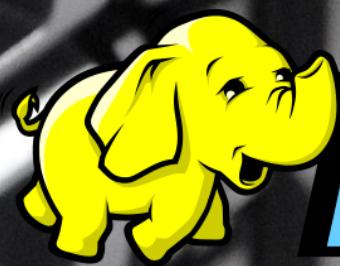




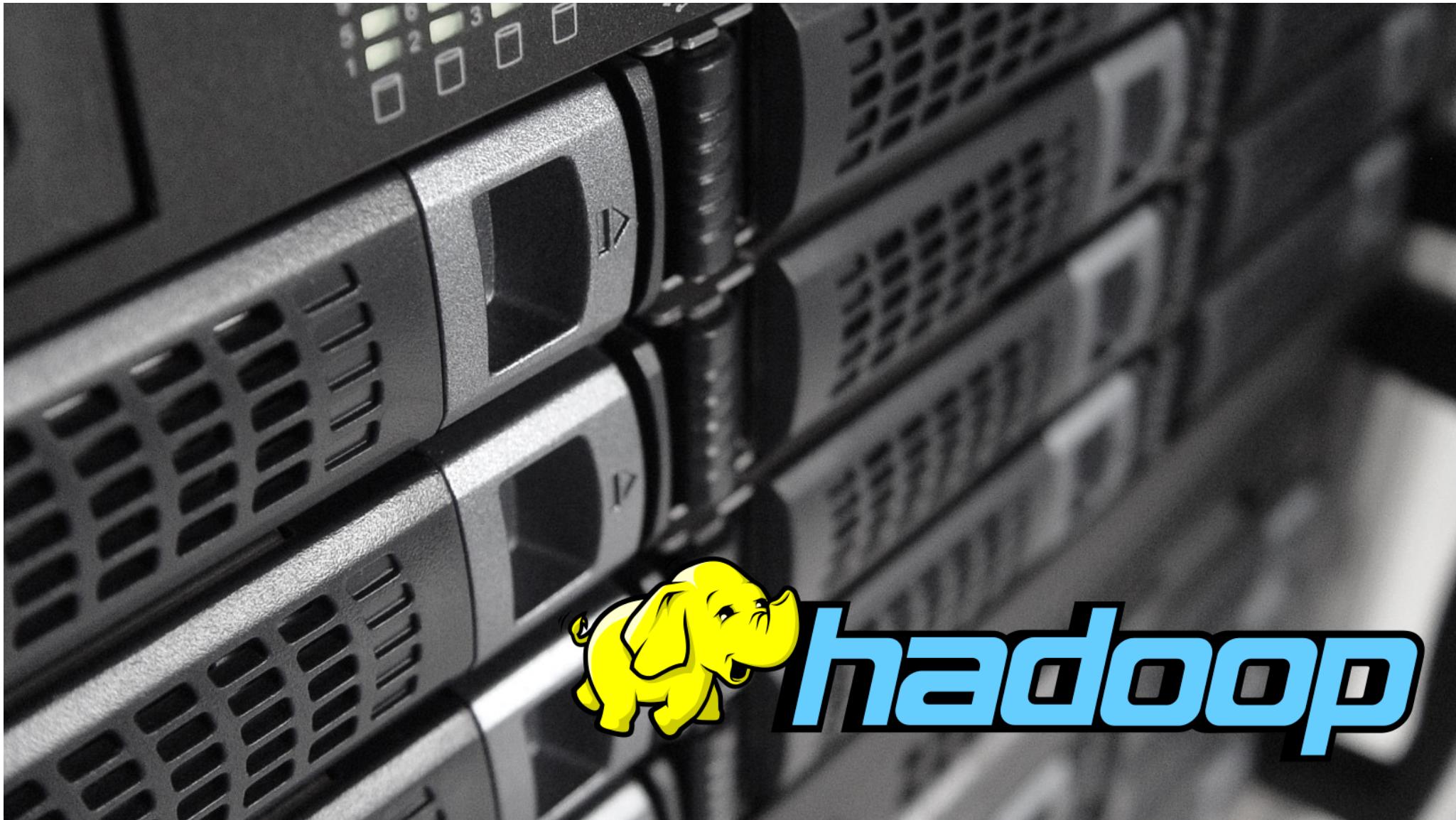




TOON®



hadoop





Requirements

- Seamless transition
- Ensure data anonymity
- Move fast, optimise later
- Ensure multi-tenancy
- As little disturbance as possible to the DS team



TOON®



TOON®



TOON®



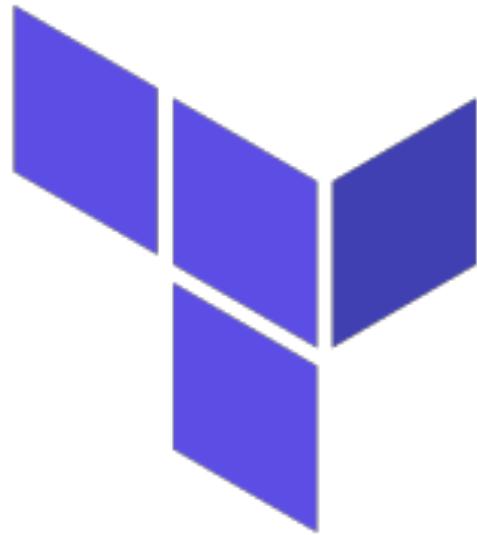
- Cluster timeouts
- Autoscaling
- Spot instances
- Well documented API

TOON®

Infrastructure as code

- Repeatability
- Fast deployment
- Resilience
- Documentation

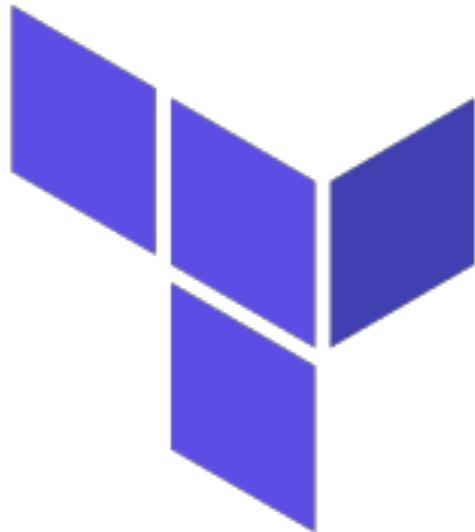




TOON®



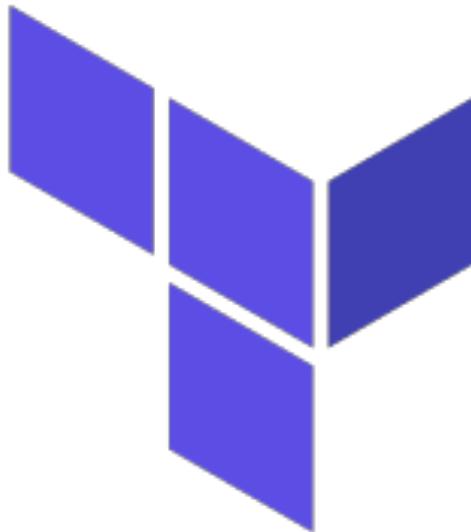
Terraform



- S3 Buckets
- EC2 instances
- Network topology
- Log management
- RDS instances
- IAM roles/policies

TOON®

Terraform



- S3 Buckets
- EC2 instances
- Network topology
- Log management
- RDS instances
- IAM roles/policies

```
./bin/provision relational-db <target> apply
```

```
./bin/provision buckets-configuration <target> apply
```

TOON®

Ansible

- User management
- Databases and ACLs
- Custom app deployment



Ansible

- User management
- Databases and ACLs
- Custom app deployment

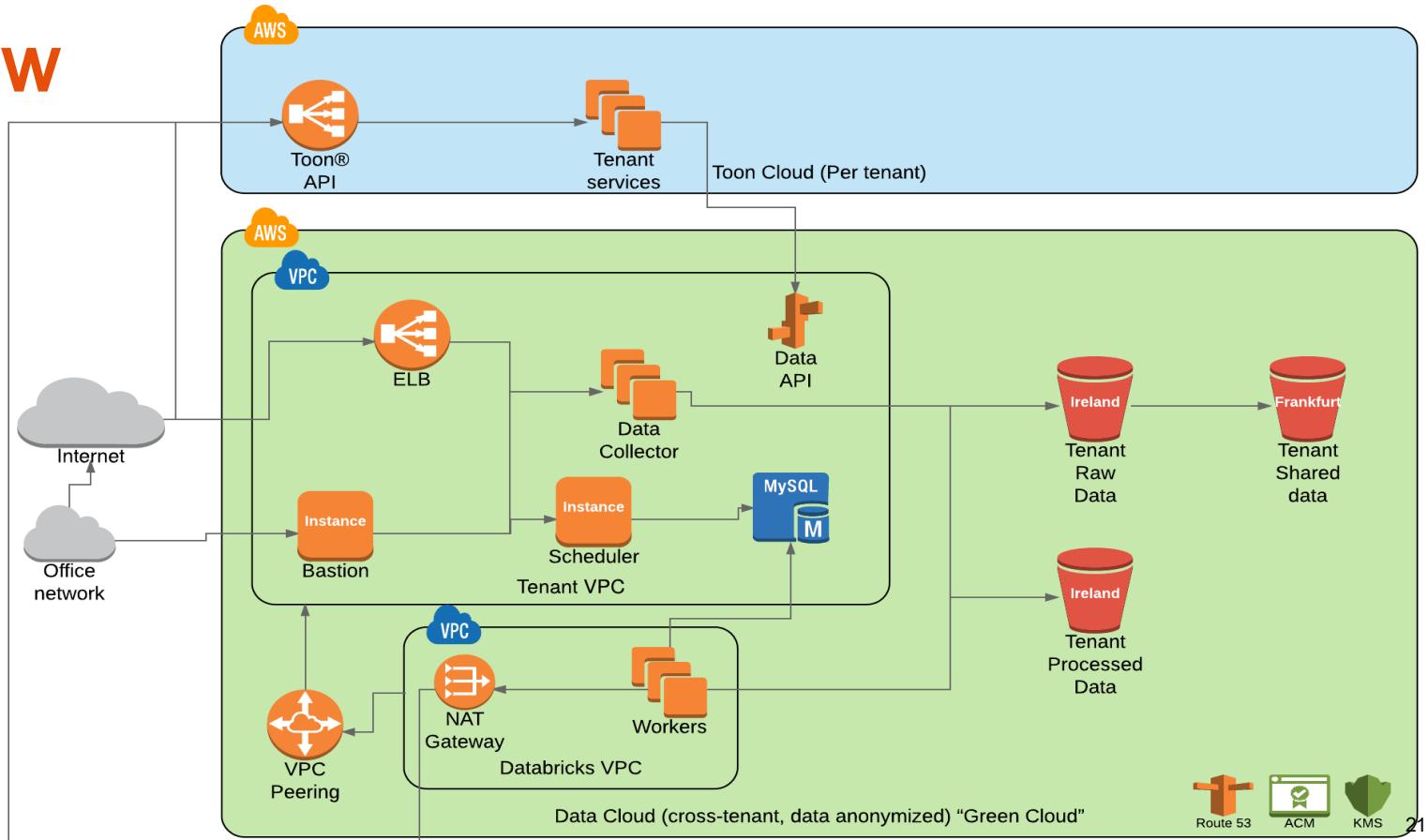


```
./bin/configure hive-metastore.yml <target>
```

```
./bin/configure create-cluster.yml <target>
```

Architecture Overview

TOON®



8 Mi

9 Do

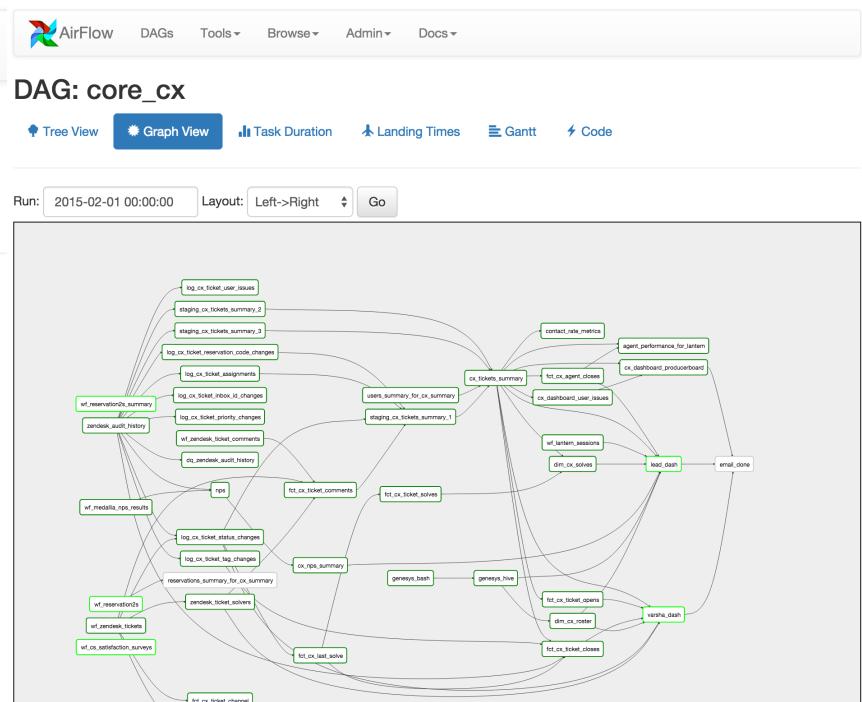
10 Fr

11 Sa

12 So

13 Mo

Airflow



TOON®

```
dag_config:  
    dag_id: <tenant>_preprocessed_dag  
    schedule_interval: '30 5 * * *'  
    start_date: 2017-08-28  
    max_active_runs: 1
```

```
task_defaults:  
    cluster_name: '<tenant>_preprocessed_dag'  
    databricks_conn_id: 'databricks_<tenant>'  
    retries: '3'  
    jar: '<location_of_jar>'
```

```
tasks:
  download_<user_group>_acc:
    class: eu.toon.etl.preprocessed.LoadPreProcessedData
    parameters:
      endpoint: <data_endpoint>
      source_db: <tenant_db>
      source_table: <table_name>
      destination_db: <tenant_db>
      destination_table: <target_table_name>
      authenticator: <authenticaton_type>
      aws_region: <aws_region>
      source: acc
      service_name: <service_name>
      aws_wso2_initial_token_key: <token_key_on_ssm>
      aws_wso2_auth_endpoint: <auth_endpoint>
      date: '{{ ds }}'
      depends_on:
        - <task_name>
```

```
new_cluster:
  spark_version: '3.2.x-scala2.11'
  node_type_id: 'i3.xlarge'
  autotermination_minutes: '20'
  spark_conf:
    spark.hadoop.java.jdo.option.ConnectionDriverName: '{{<hive_driver>}}'
    spark.hadoop.java.jdo.option.ConnectionUserName: '{{<hive_username>}}'
    spark.hadoop.java.jdo.option.ConnectionPassword: '{{<hive_passowrd>}}'
    spark.hadoop.java.jdo.option.ConnectionURL: '{{<hive_url>}}'
    spark.serializer: 'org.apache.spark.serializer.KryoSerializer'
    spark.hadoop.fs.s3.impl: 'com.databricks.s3a.S3AFileSystem'
    spark.sql.hive.metastore.jars: 'builtin'
    spark.io.compression.codec: 'snappy'
    spark.sql.hive.metastore.version: '1.2.1'
  aws_attributes:
    first_on_demand: '2'
    availability: 'SPOT'
    zone_id: '<aws_region>'
    instance_profile_arn: '<arn_of_profile>'
    spot_bid_price_percent: '75'
  num_workers: '5'
  custom_tags:
    Environment : '<tenant>.prod'
    GroupName : 'BigData'
    Tenant: '<tenant_name>'
    Provisioner : 'Databricks'
    CostCenter : 'international-operations'
    Stage : 'prod'
```



TOON[®]



- External Hive metastore
- Send logs to S3
- Authorisation
- i3.2xlarge nodes

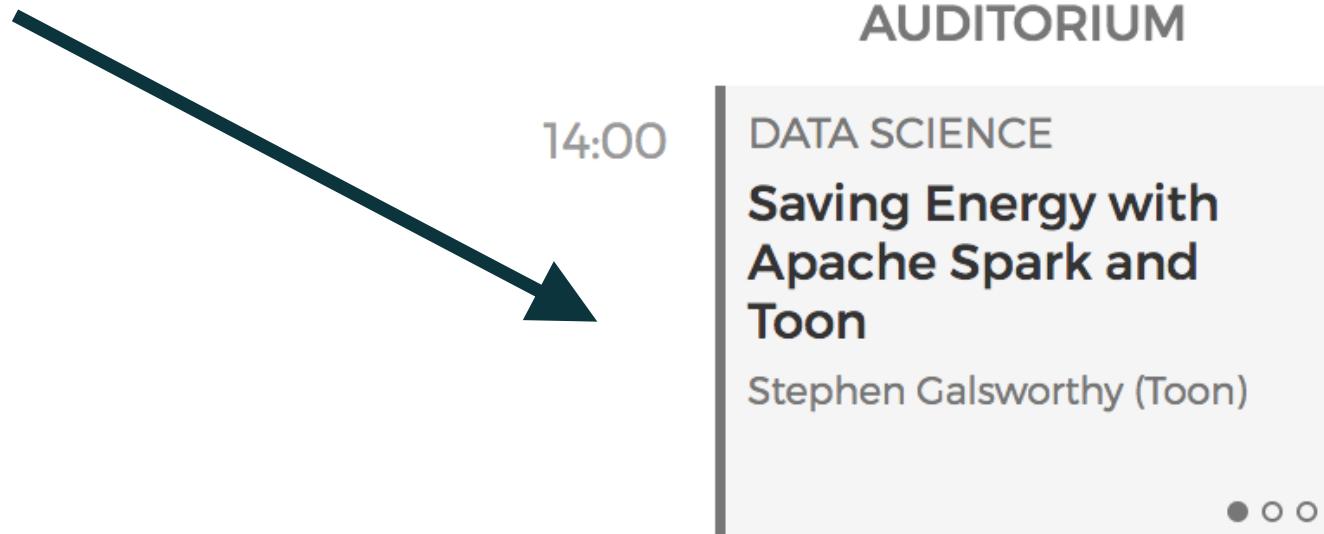
Future plans

- Streaming
- Real time services
- Improve CI/CD

What's all this for?

TOON[®]

What's all this for?



TOON®

31

Thanks to the team

Aemro Amare

Barend Garvelink

Bert Jan Katsman

Kliment Markovski

Miquel Monreal

Stanislava Potupchik



TOON®

Questions?

TOON®

