

# Building Smart IoT Applications with Python and Spark

...

Rafael Schultze-Kraft  
PyData Berlin 2017



# WHO WE ARE



WATT *x*

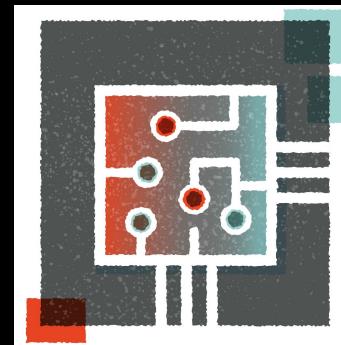


# OUR PROCESS

Ideation



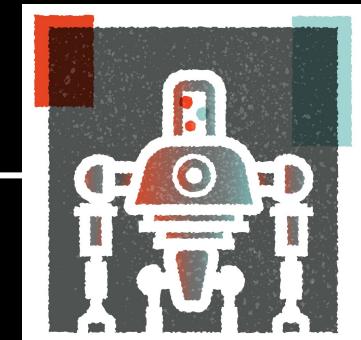
Prototyping



Discovery



Venture



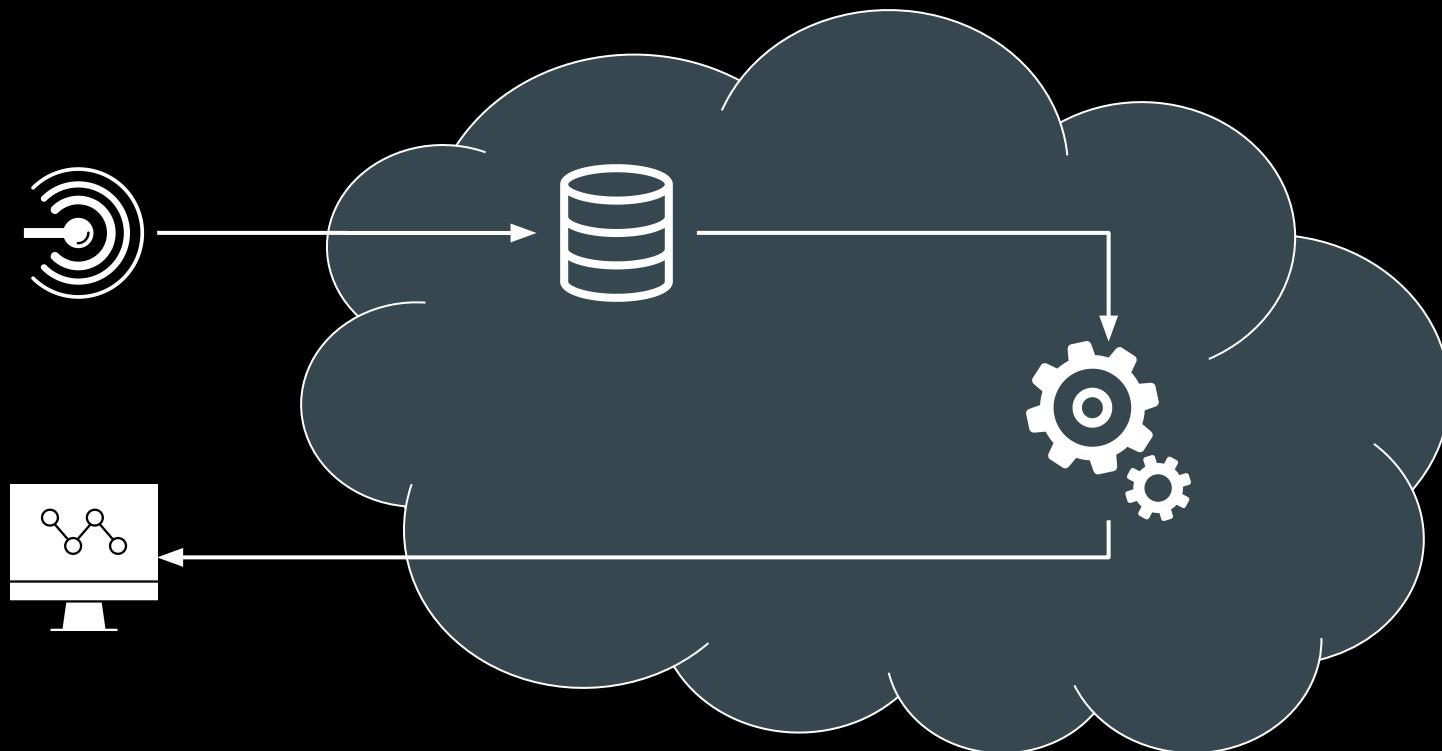
# Smart Building Infrastructure



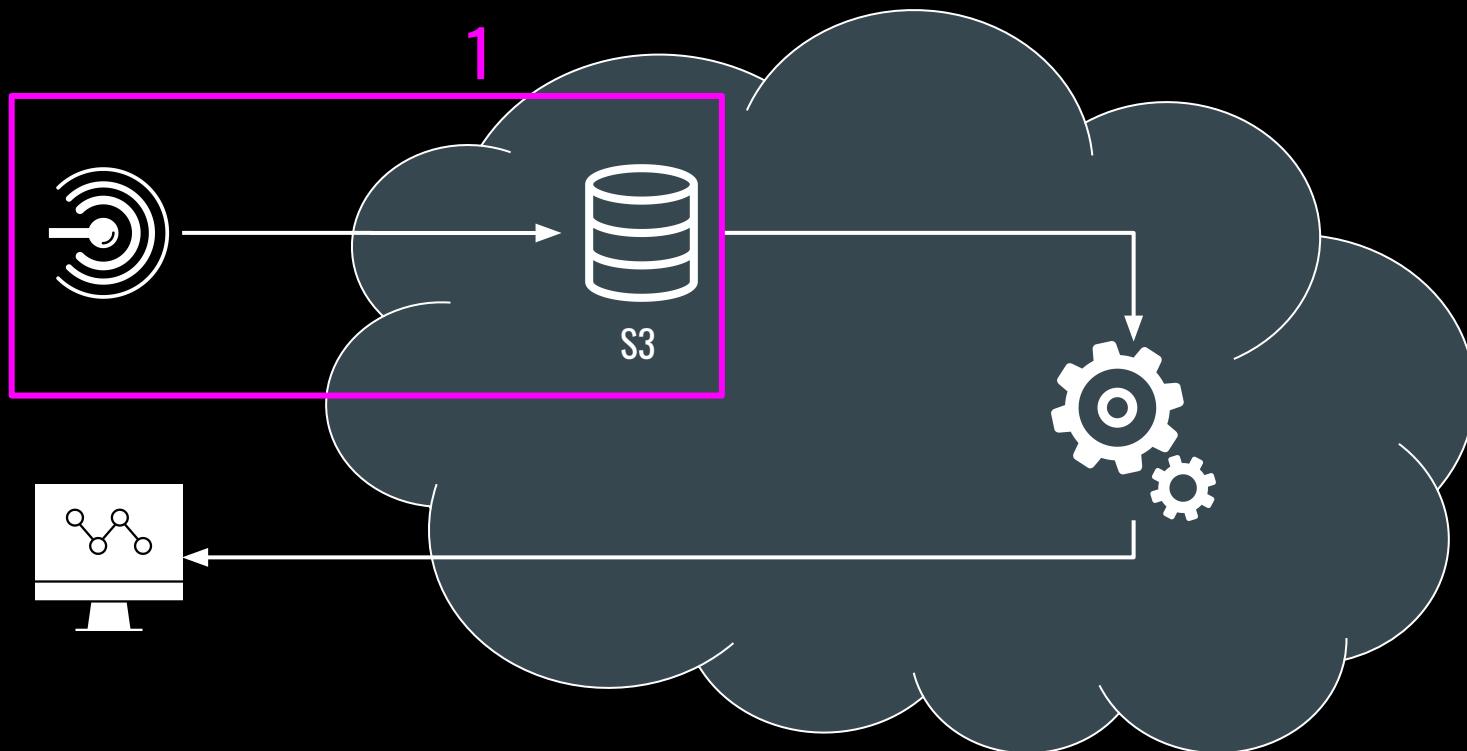
# SNUK

[www.snuk.io](http://www.snuk.io)

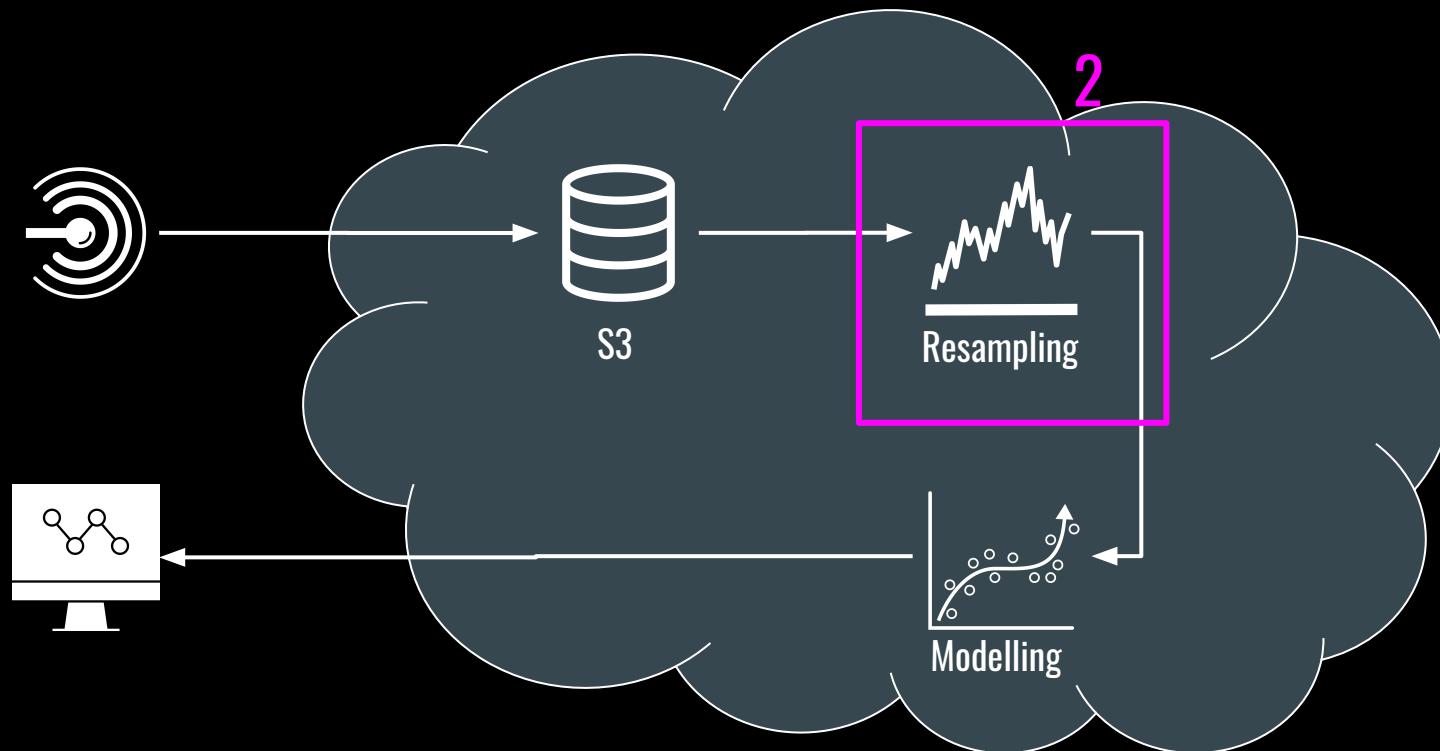
# IoT APPLICATION



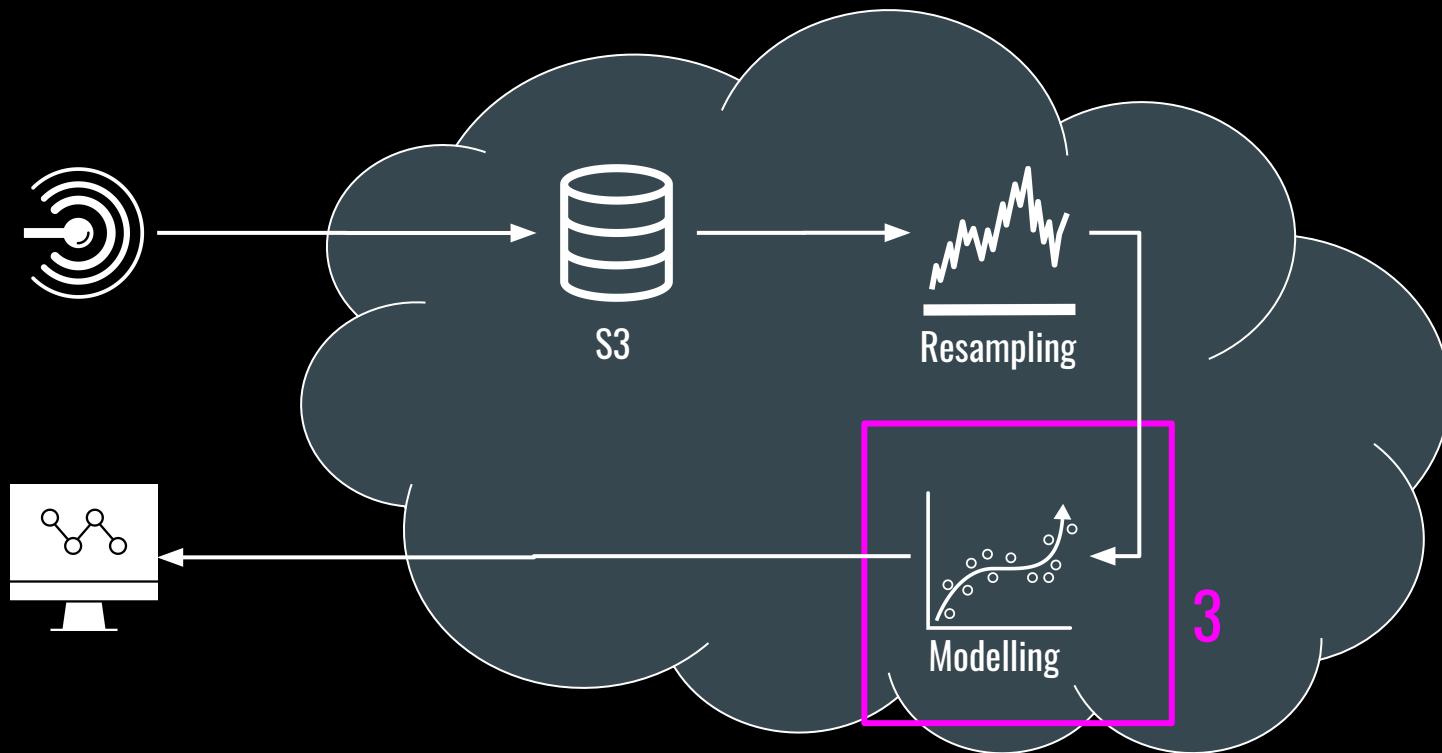
# 1 SMART OFFICE IoT SETUP



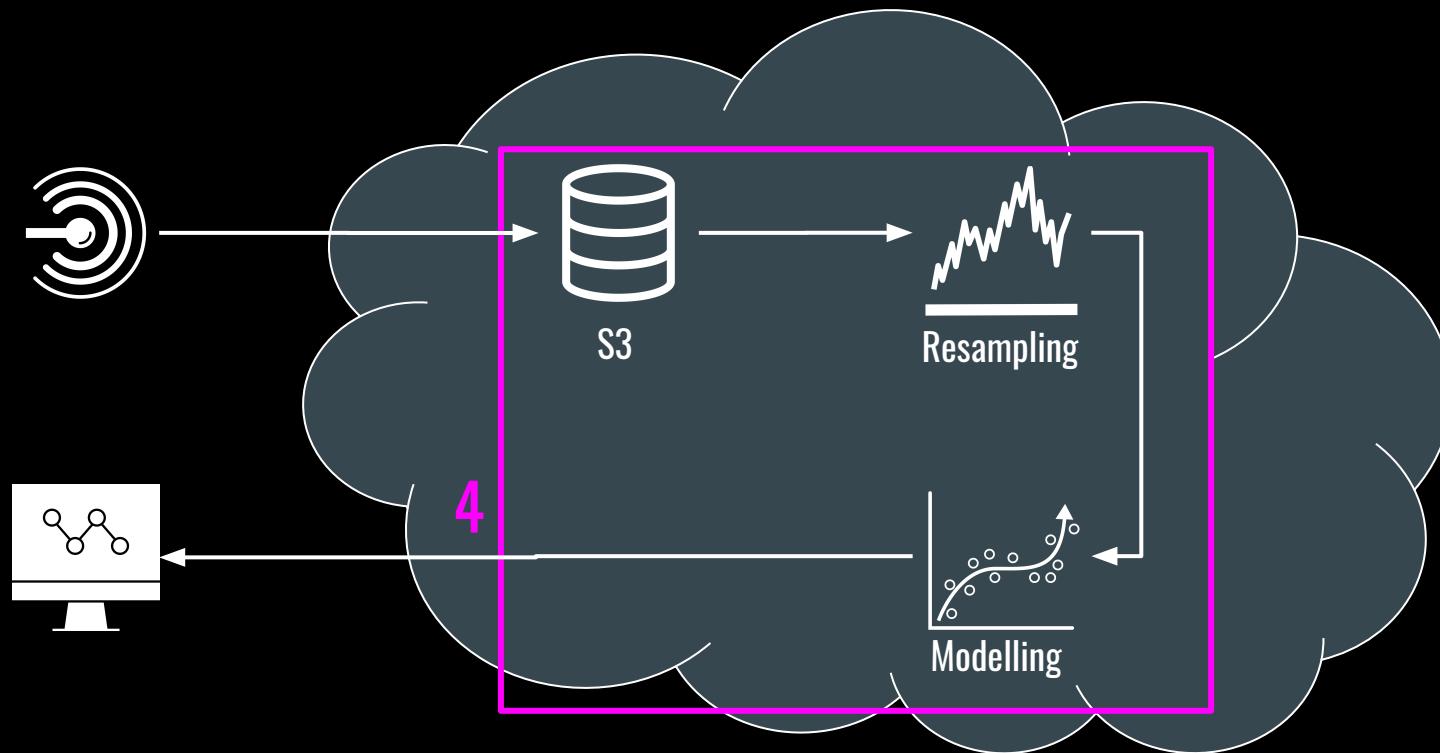
# 2 RESAMPLING IoT SENSOR DATA



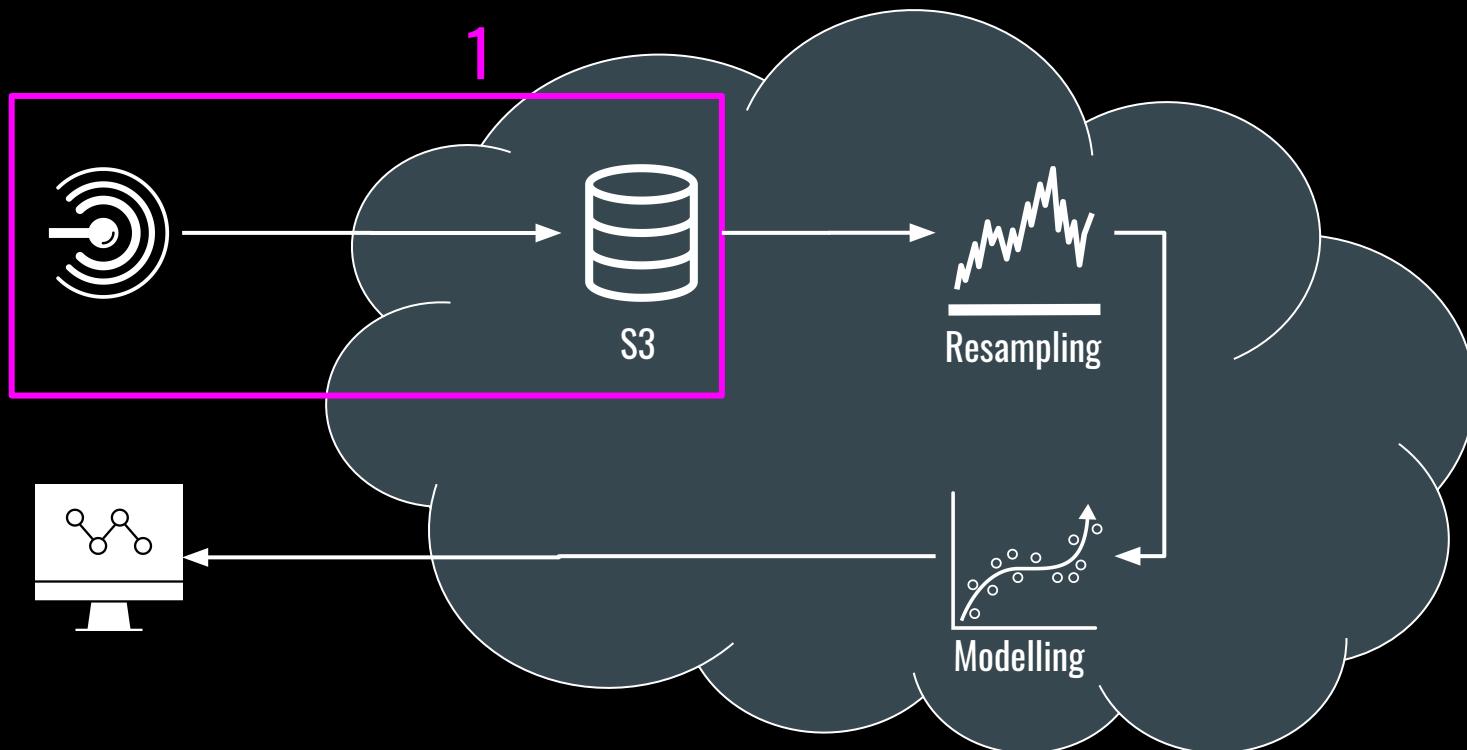
# 3 MODELLING IoT SENSOR DATA



# 4 SCALING AND AUTOMATING



# 1 SMART OFFICE IoT SETUP



# SETUP



4000m<sup>2</sup> office space



2 floors



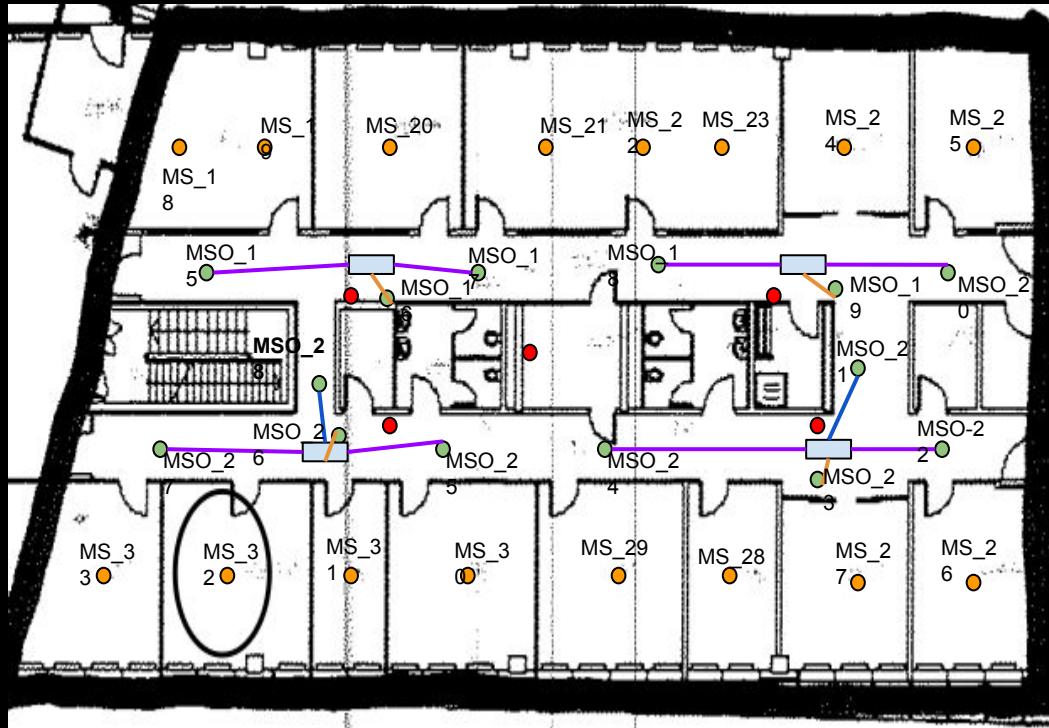
36 rooms



215 sensors



3 months of data



# SENSORS



Temperature



Luminosity

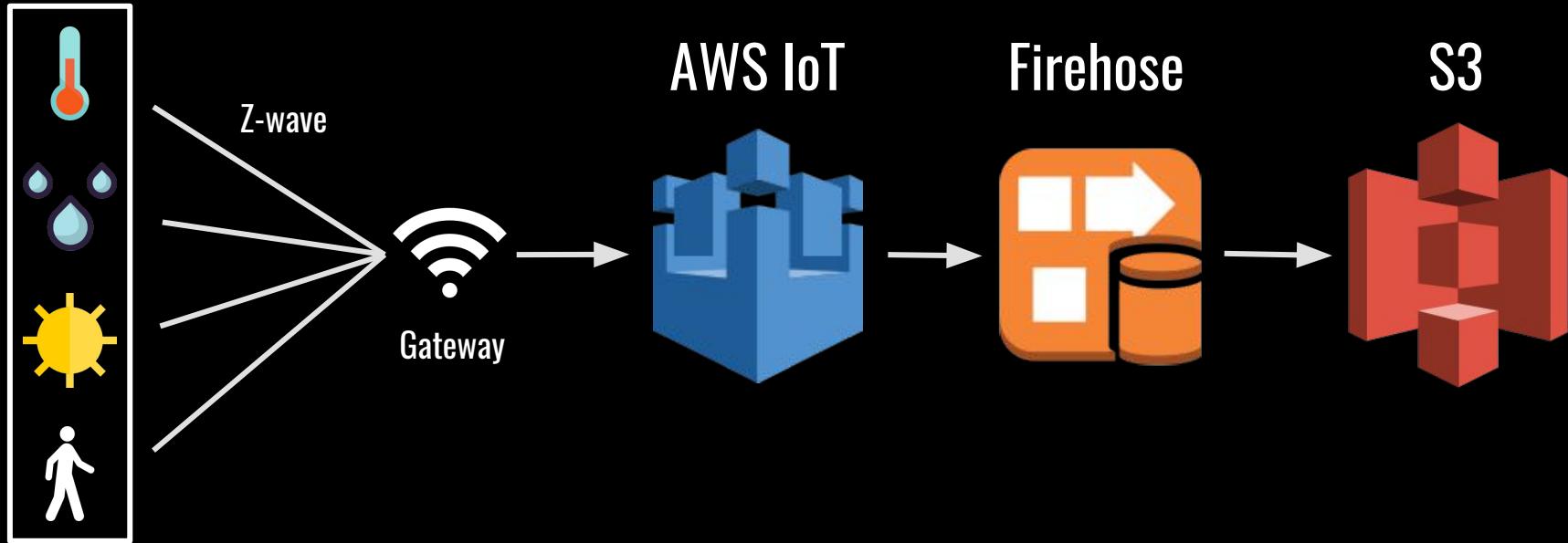


Humidity



Motion

# DATA FLOW



# RAW DATA

s3://raw\_data/2016/08/14/

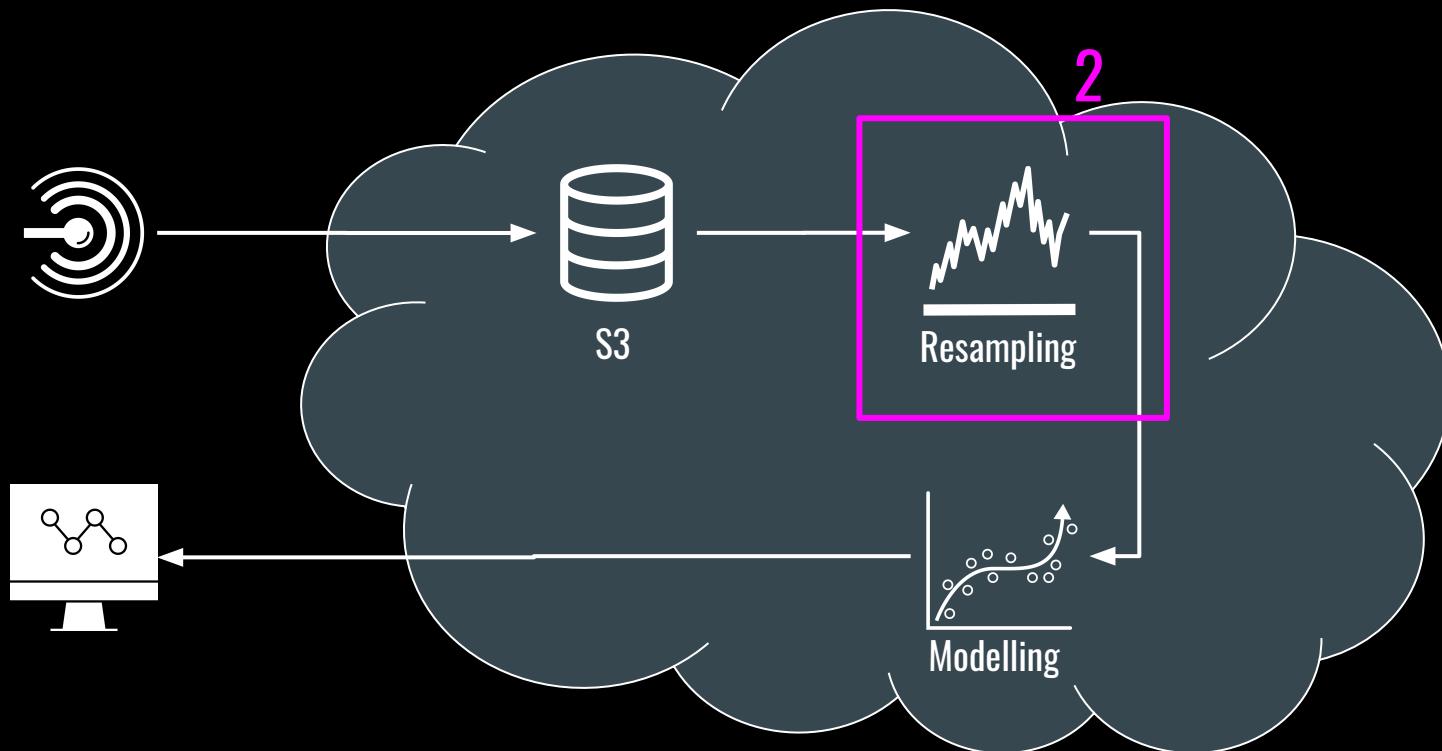
```
sensor_id='8514-294-220-250-m', timestamp=1471020341000, value=1.0,  
sensor_id='7513-133-195-250-m', timestamp=1471020345000, value=0.0,  
sensor_id='8514-294-220-250-t', timestamp=1471132284000, value=25.0,  
sensor_id='7513-133-195-250-t', timestamp=1471132306000, value=23.9,  
sensor_id='7513-133-195-250-t', timestamp=1471132316000, value=23.0,  
sensor_id='7505-277-172-250-m', timestamp=1471132385000, value=0.0,  
sensor_id='7573-230-160-250-l', timestamp=1471132448000, value=4.0,
```

.

.

.

# 2 RESAMPLING IoT SENSOR DATA



# OBSERVATIONS to TIME SERIES

`sensor_id, timestamp, value`

•           •           •

•           •           •

•           •           •



`timestamp, sensor_1, sensor_2, sensor_3, ...`

$t_1$            •           •           •           •

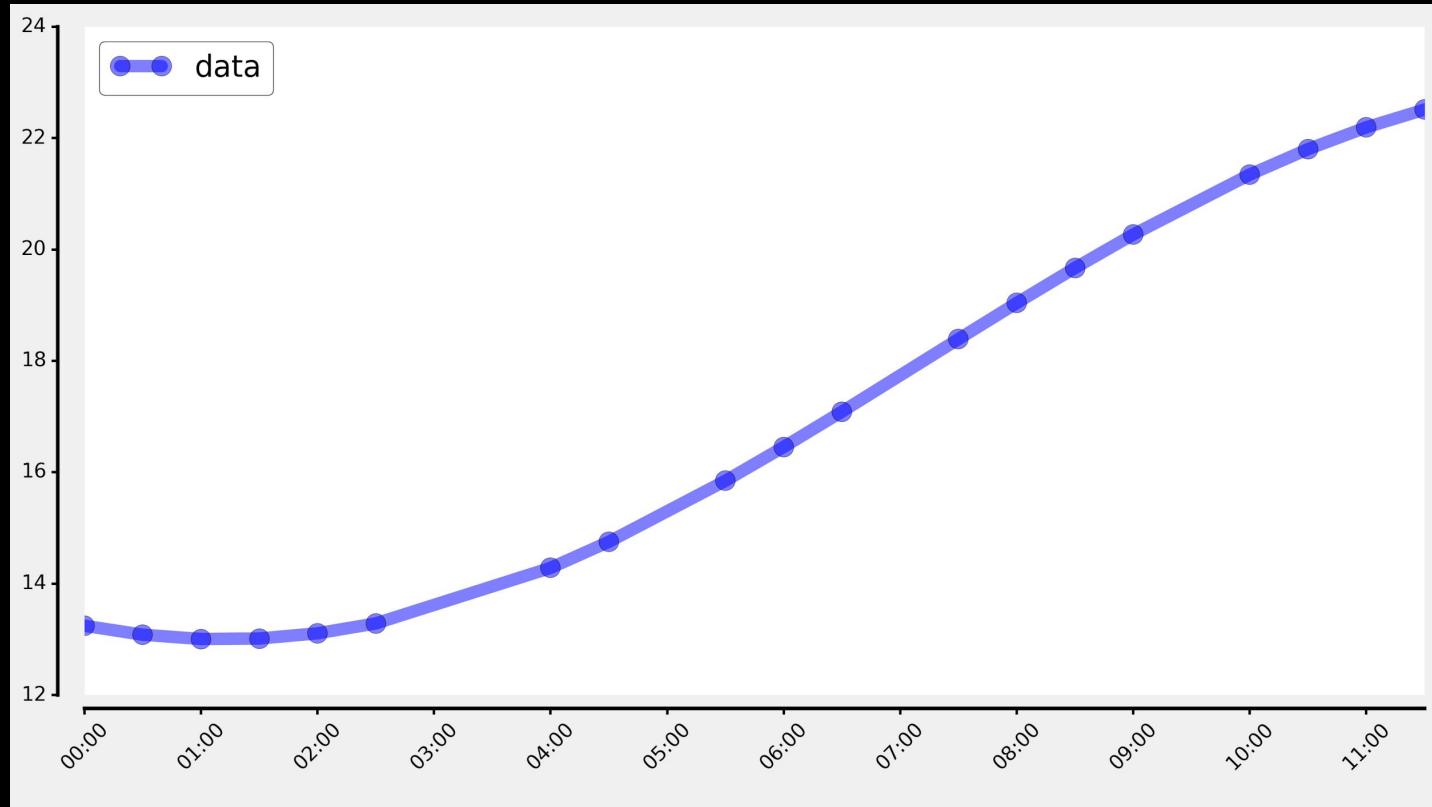
$t_2$            •           •           •           •

$t_3$            •           •           •           •

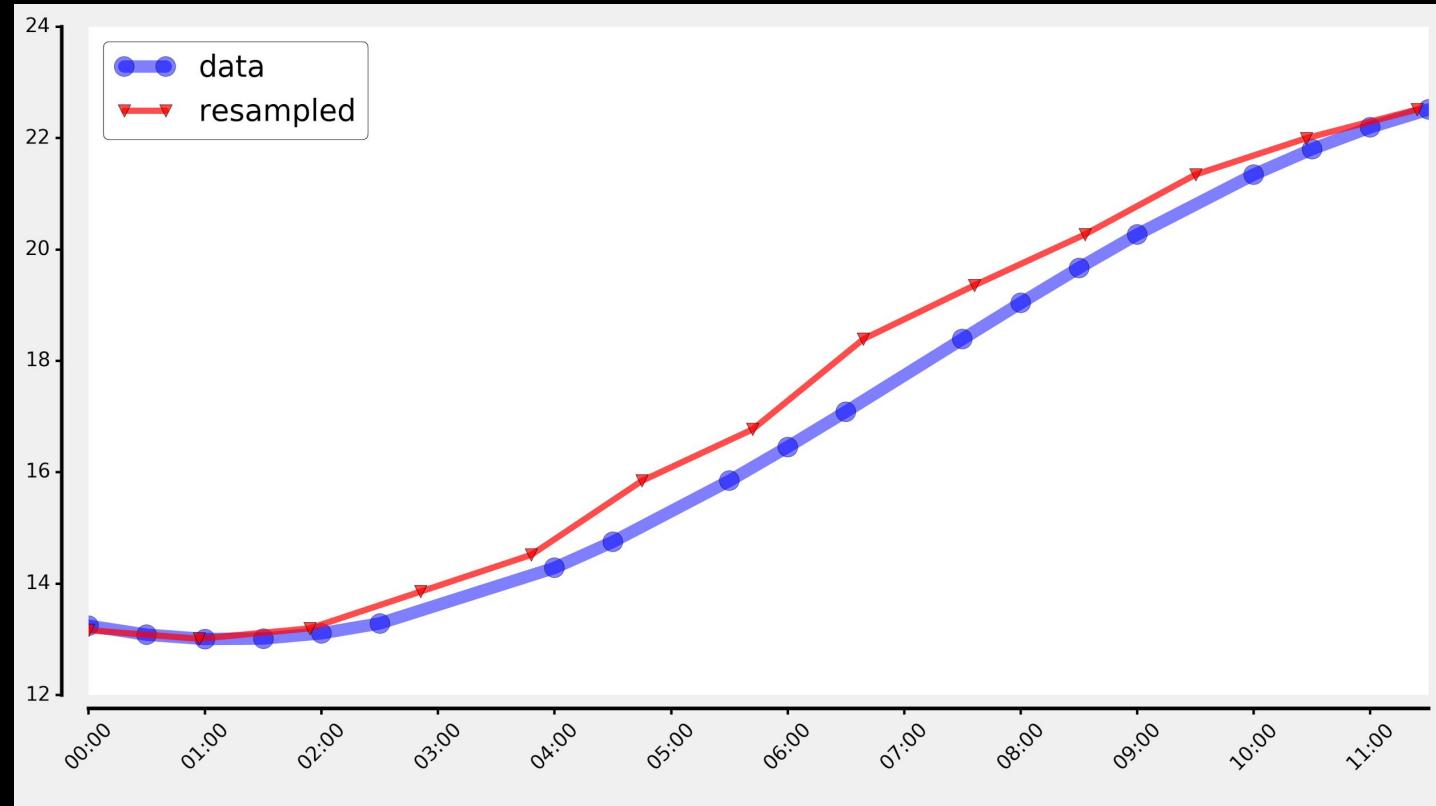
# RESAMPLE? PANDAS!

```
import pandas as pd  
  
df = df.resample(df, '5min', 'mean')
```

# RESAMPLING



# RESAMPLING

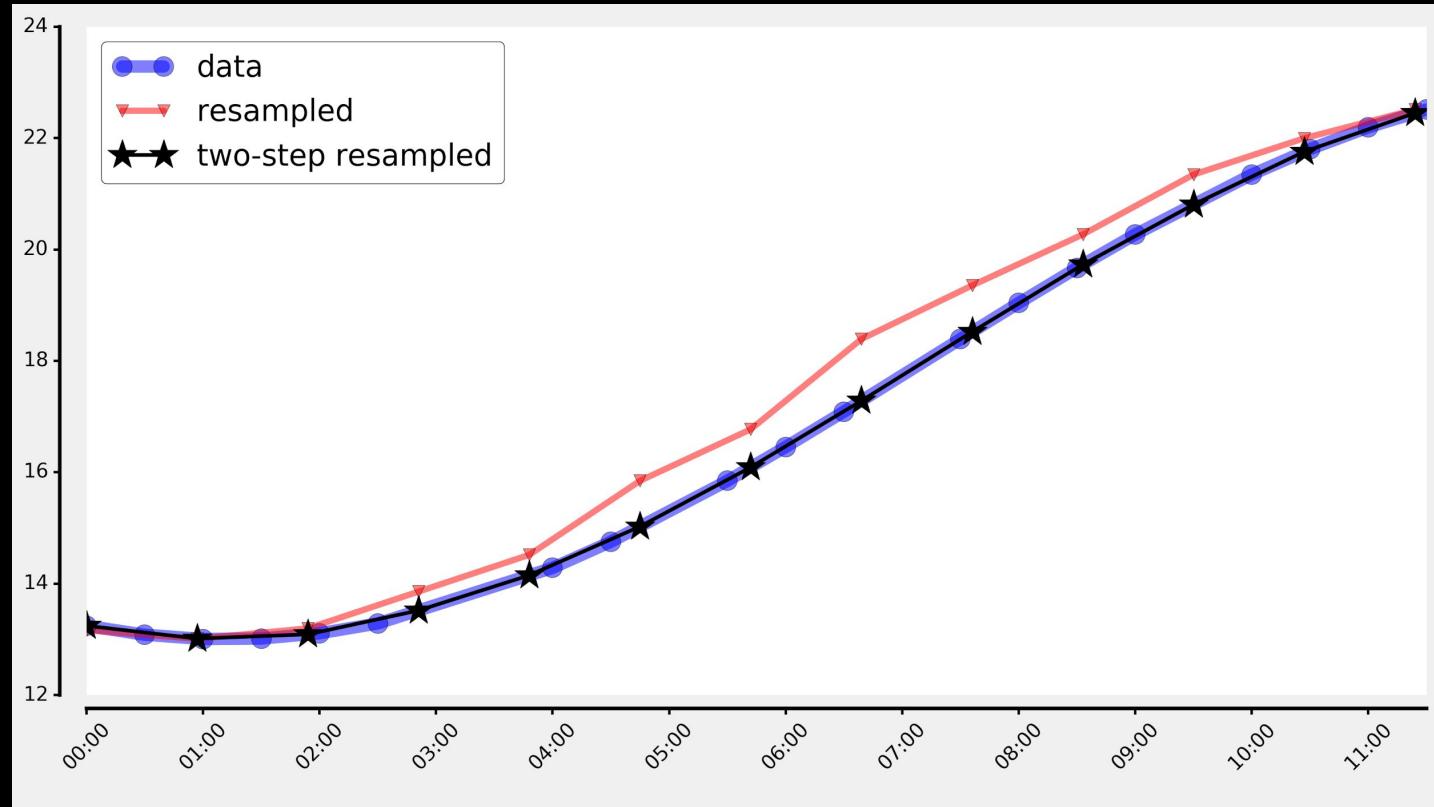


# TWO-STEP RESAMPLING PROCESS

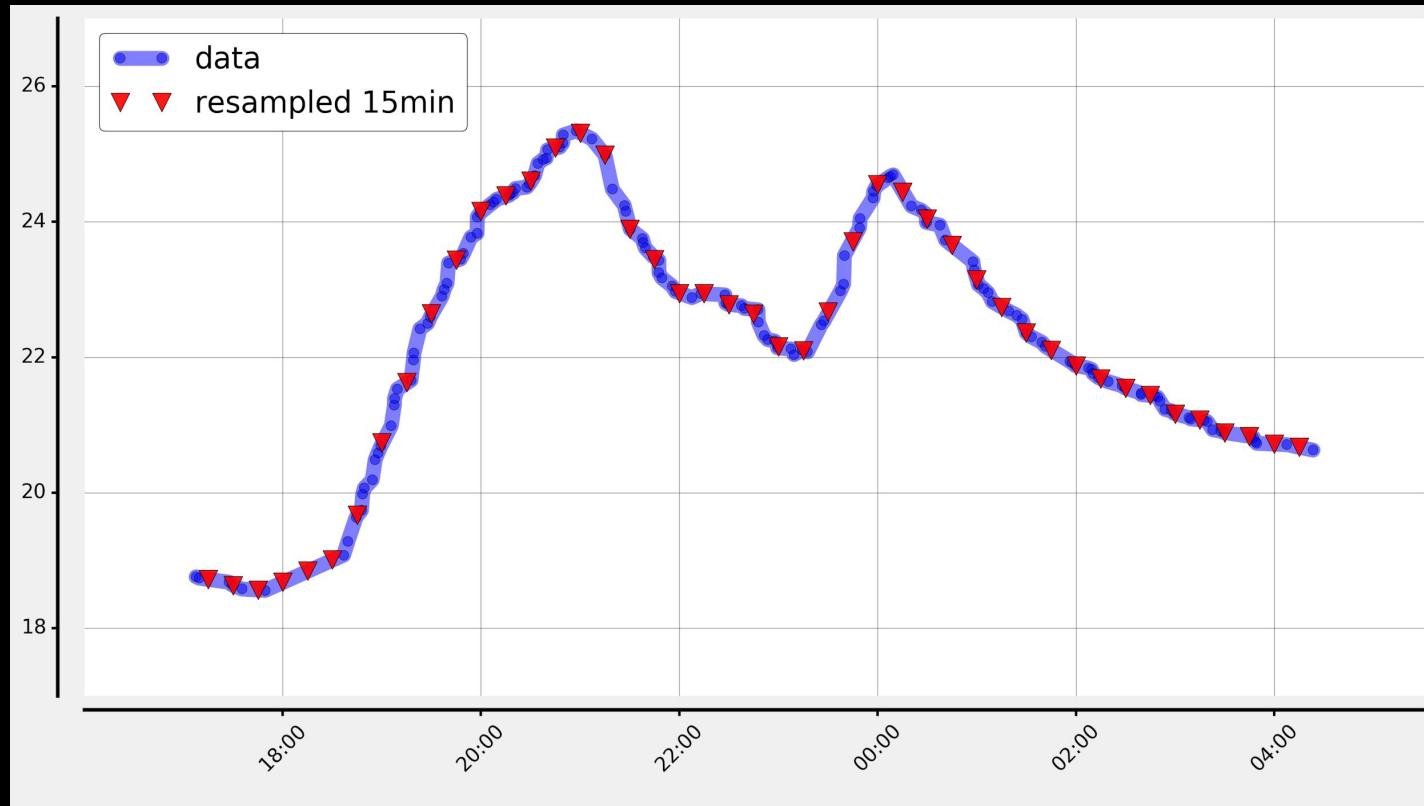
```
# 1) upsample to high frequency
tmp = data.resample('1min', 'mean').interpolate()

# 2) downsample to desired frequency
resampled = tmp.resample('15min', 'ffill')
```

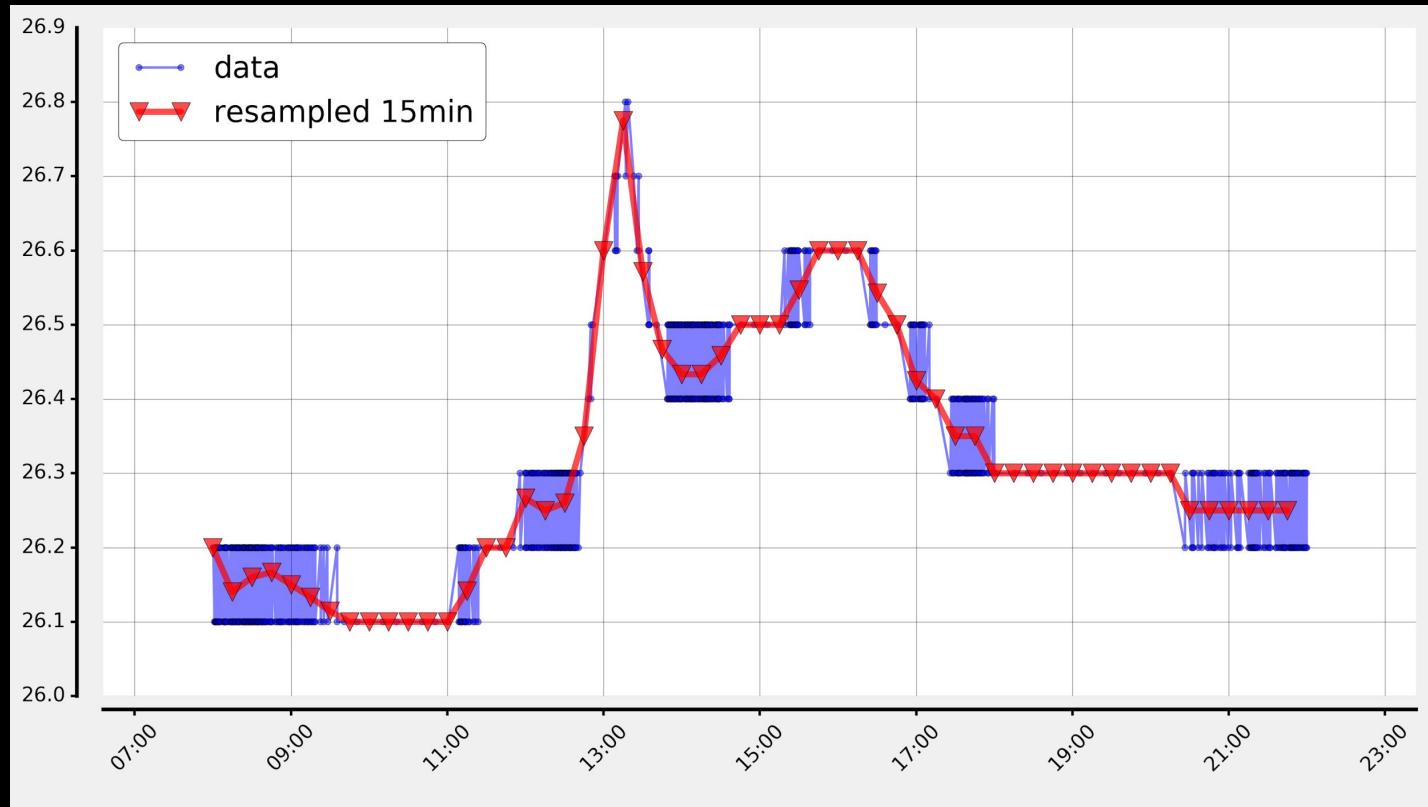
# RESAMPLING



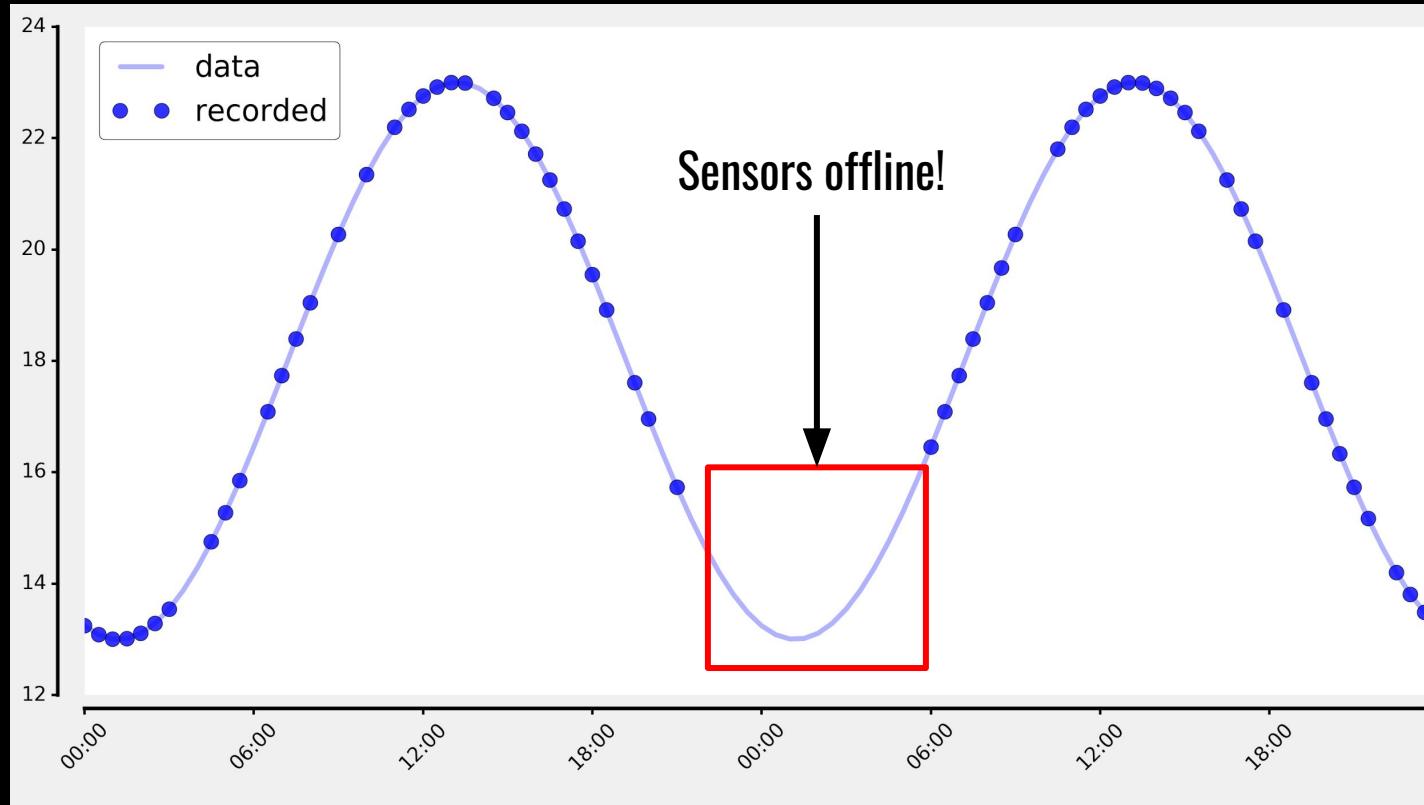
# RESAMPLING



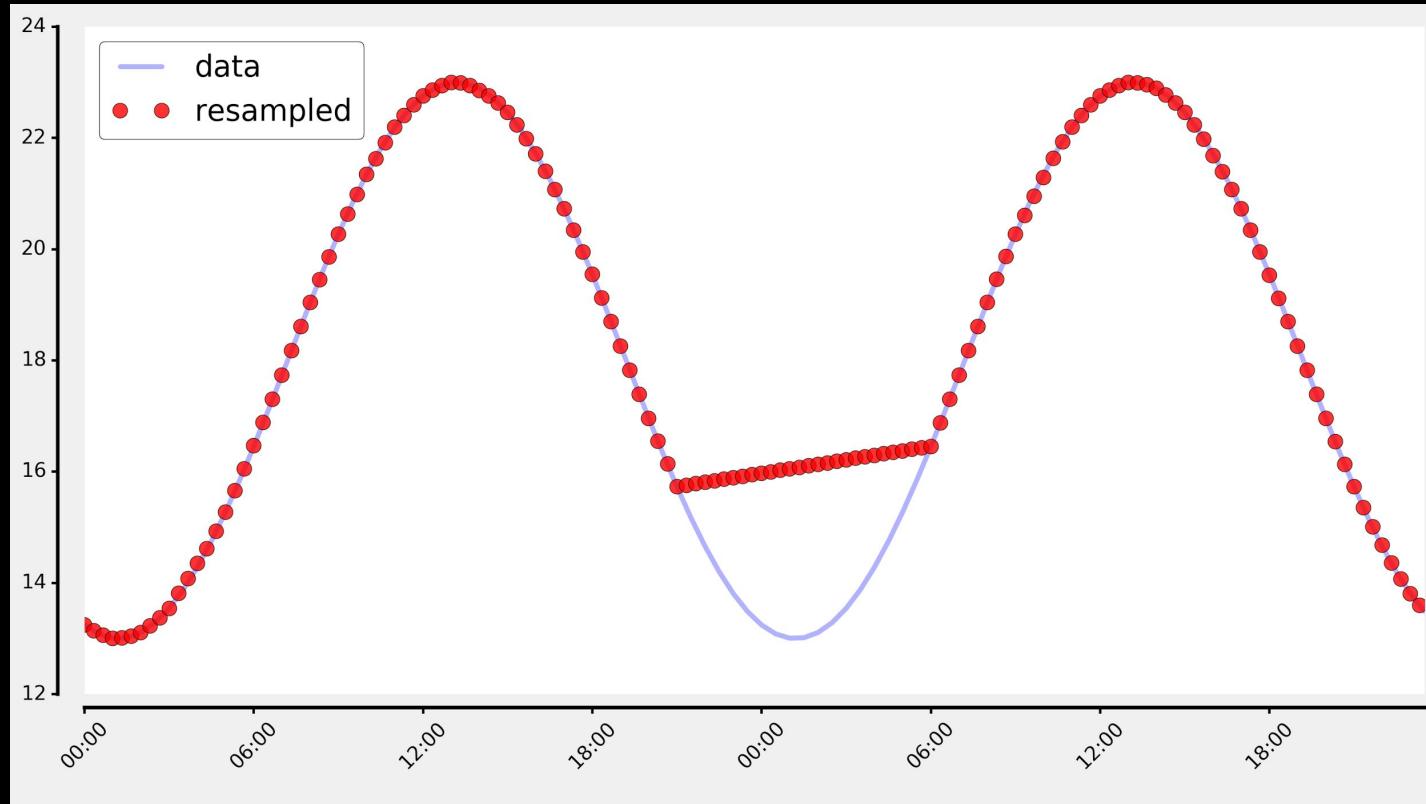
# RESAMPLING



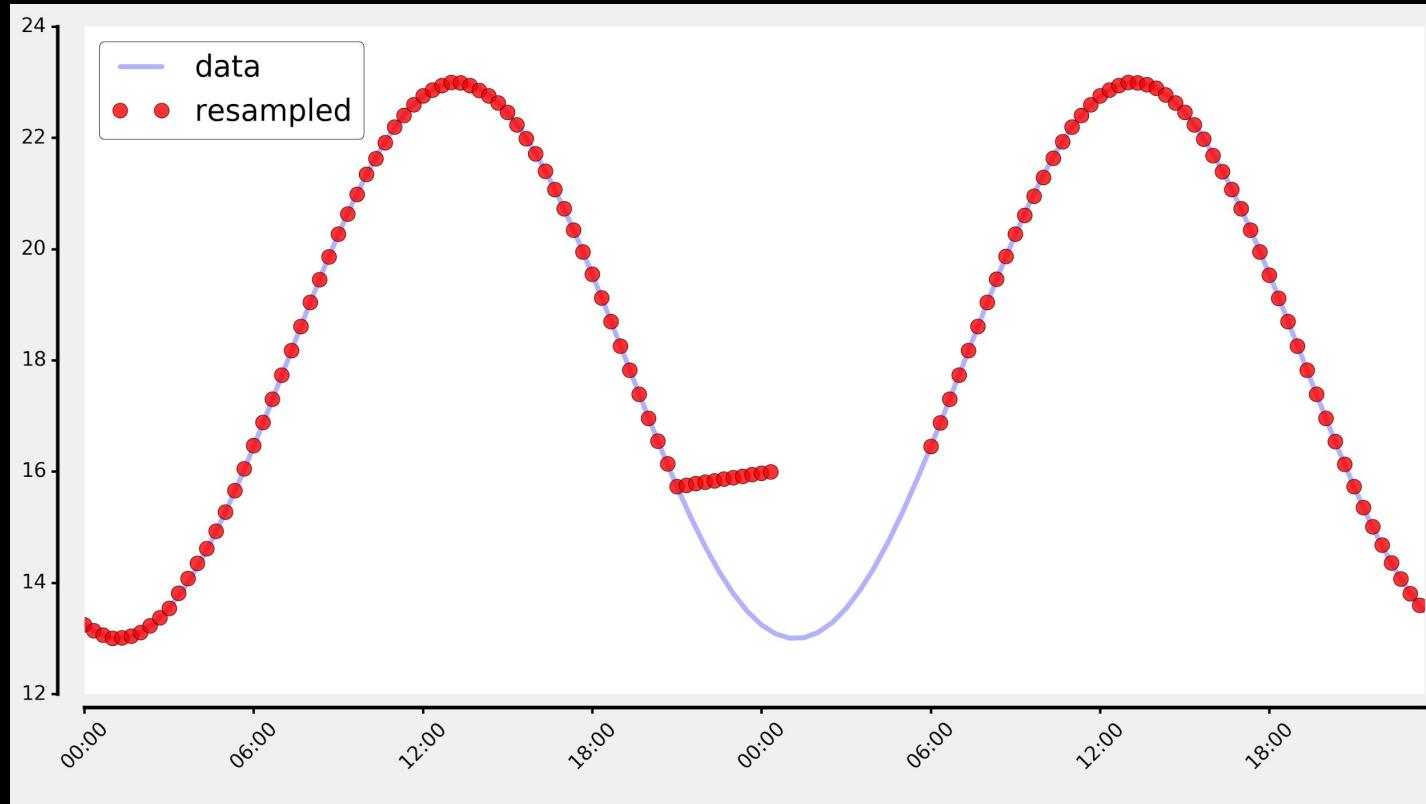
# PROBLEM: OFFLINE SENSORS



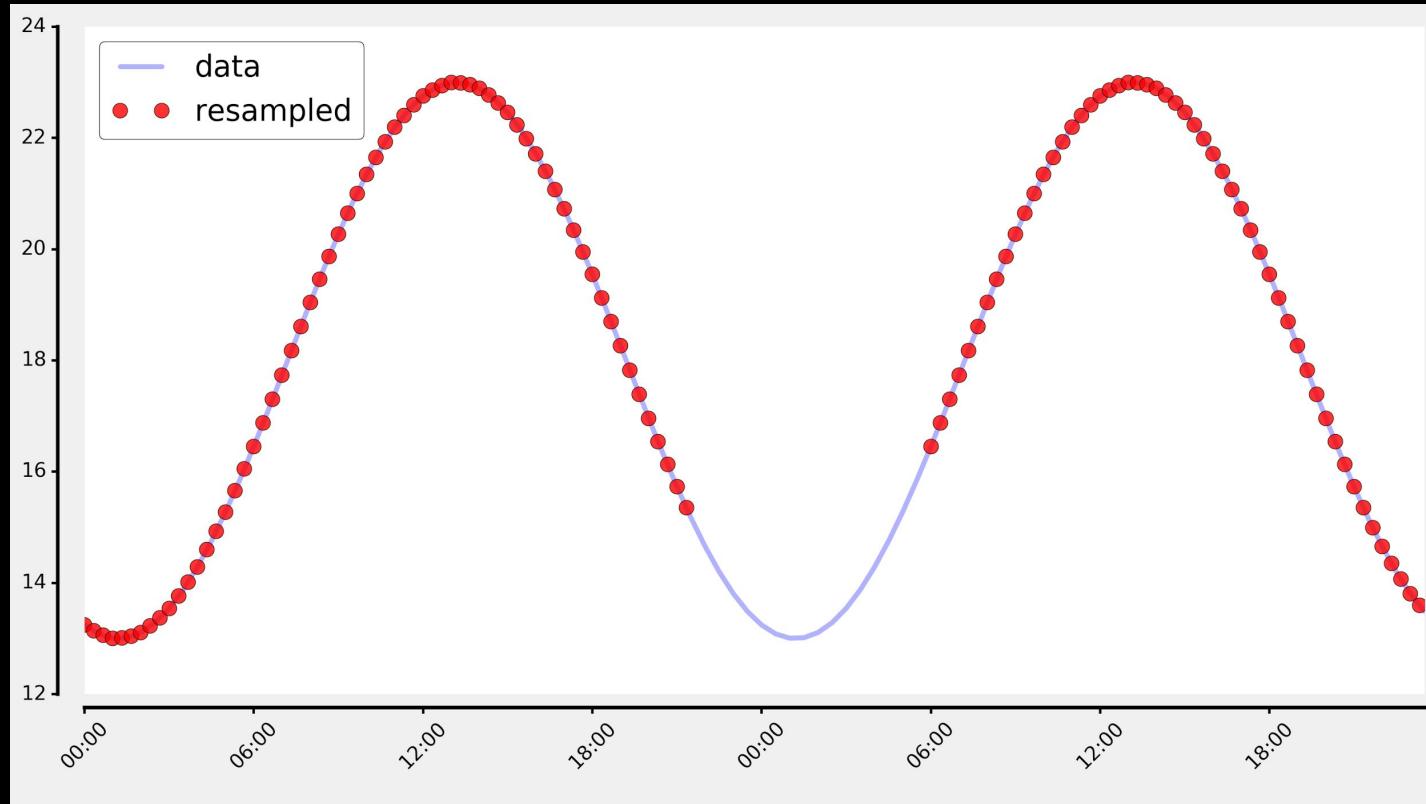
# OFFLINE SENSORS



# OFFLINE SENSORS



# OFFLINE SENSORS



# TIME SERIES

**2016/06/20**

**occ-7115**

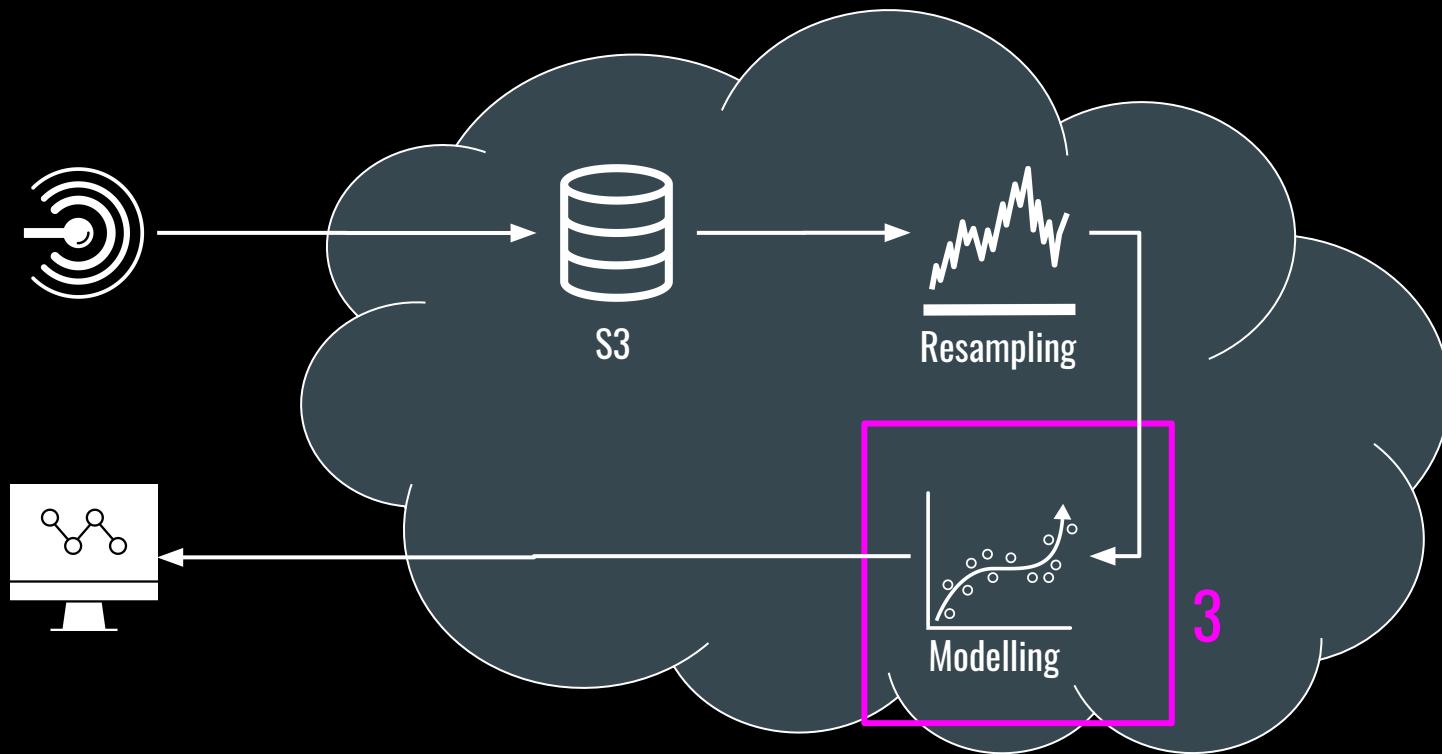
2016-06-20	11:00:00	0.000000
2016-06-20	11:05:00	0.000000
2016-06-20	11:10:00	0.056667
2016-06-20	11:15:00	0.770000
2016-06-20	11:20:00	0.766667
2016-06-20	11:25:00	0.880000
2016-06-20	11:30:00	0.060000
2016-06-20	11:35:00	0.446667
2016-06-20	11:40:00	0.790000
2016-06-20	11:45:00	0.860000
2016-06-20	11:50:00	0.800000
2016-06-20	11:55:00	0.440000
2016-06-20	12:00:00	1.000000

# TIME SERIES

**2016/06/20**

		<b>occ-7115</b>	<b>tmp-7115</b>
2016-06-20	11:00:00	0.000000	23.100000
2016-06-20	11:05:00	0.000000	23.120000
2016-06-20	11:10:00	0.056667	23.566000
2016-06-20	11:15:00	0.770000	23.800000
2016-06-20	11:20:00	0.766667	23.810000
2016-06-20	11:25:00	0.880000	23.810000
2016-06-20	11:30:00	0.060000	23.810000
2016-06-20	11:35:00	0.446667	23.950000
2016-06-20	11:40:00	0.790000	23.950000
2016-06-20	11:45:00	0.860000	23.955500
2016-06-20	11:50:00	0.800000	23.960000
2016-06-20	11:55:00	0.440000	23.980000
2016-06-20	12:00:00	1.000000	24.111200

# 3 MODELLING IoT SENSOR DATA

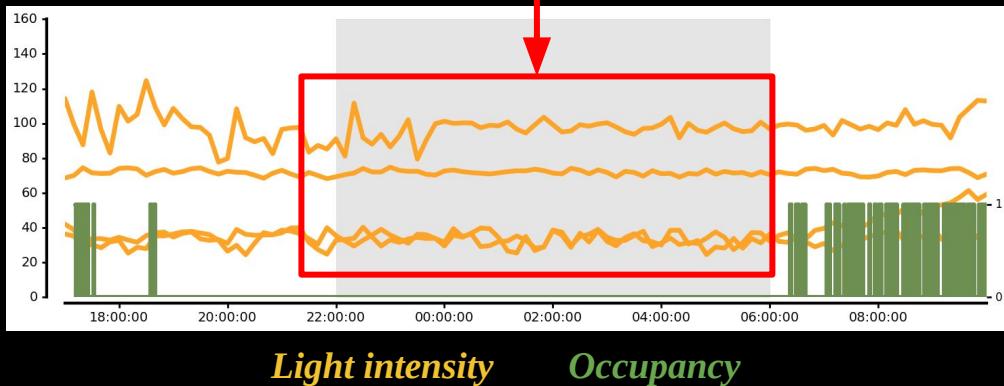


# PREDICTING OCCUPANCY

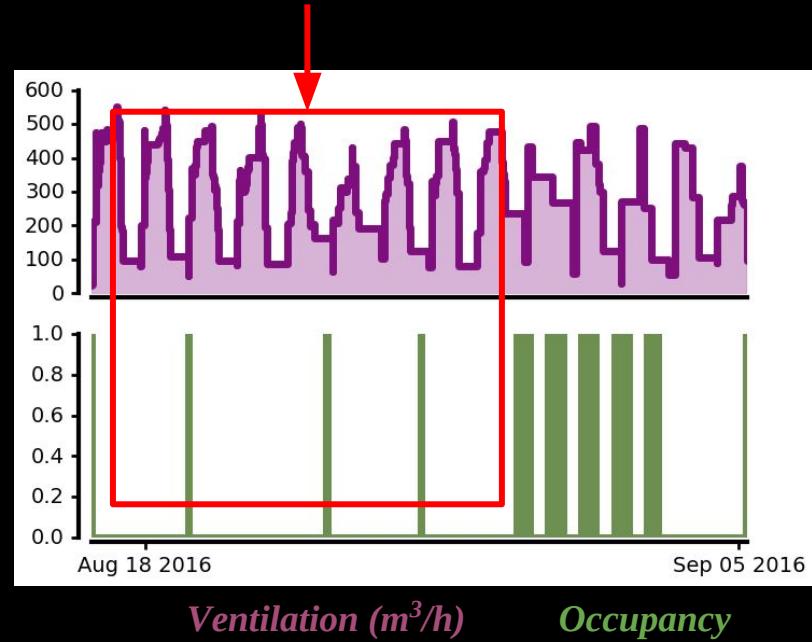


# MOTIVATION

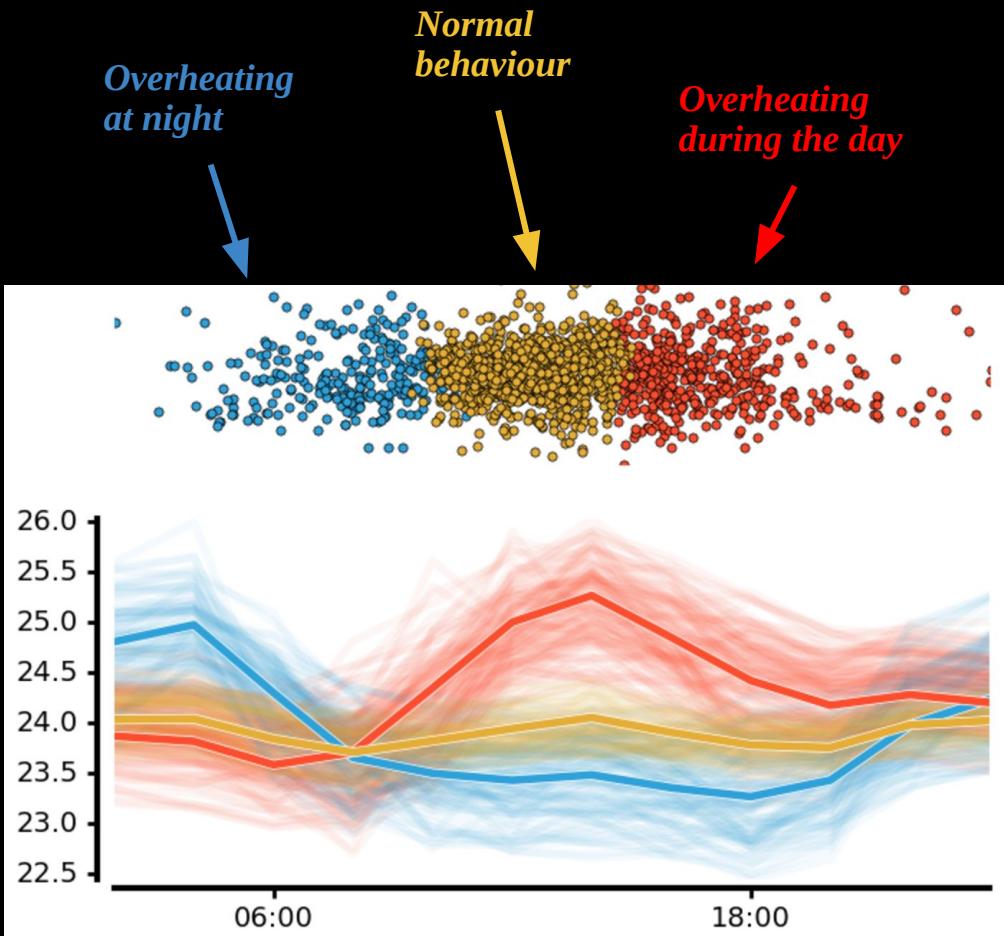
*Lights are left on at night  
when nobody is present*



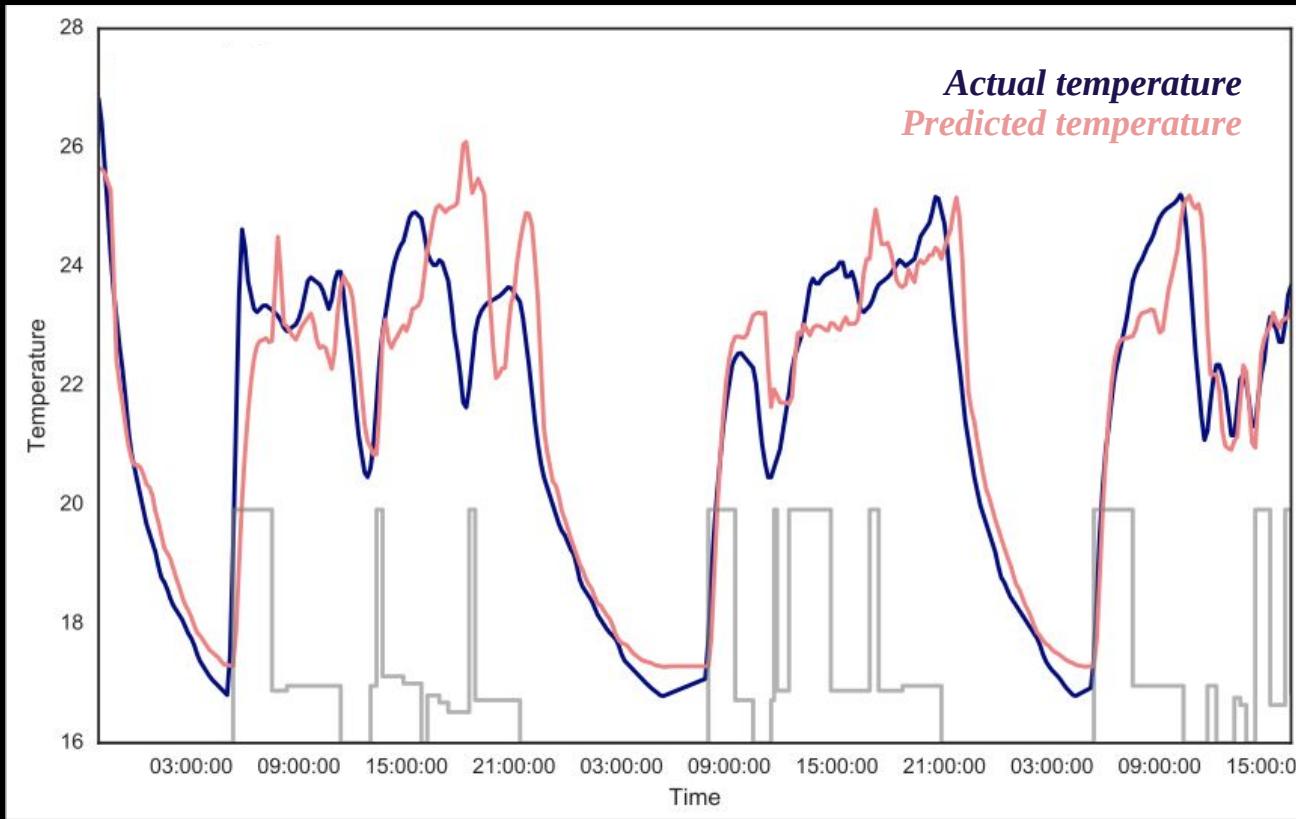
*Rooms are being cooled at full  
intensity when nobody is present*



# MOTIVATION



# MOTIVATION



# MULTI-LABEL PREDICTION



$$[x_{t+s}, \dots, x_{t+1}] = f(x_t, \dots, x_{t-n+1}) + \epsilon$$

# TIME-SERIES TO SUPERVISED LEARNING

```
import pandas as pd
```

```
df = DataFrame()  
df['t'] = range(10)  
df['t+1'] = df.t.shift(-1)  
df['t+2'] = df.t.shift(-2)  
df['t-1'] = df.t.shift(1)  
df['t-2'] = df.t.shift(2)
```

	t-2	t-1	t	t+1	t+2
	nan	nan	0	1	2
	nan	0	1	2	3
	0	1	2	3	4
	1	2	3	X	5
	2	3	4	5	6
	3	4	5	6	7
	4	5	6	7	8
	5	6	7	8	9
	6	7	8	9	nan
	7	8	9	nan	nan

# TIME-SERIES TO SUPERVISED LEARNING

```
import pandas as pd
```

```
df = DataFrame()  
df['t'] = range(10)  
df['t+1'] = df.t.shift(-1)  
df['t+2'] = df.t.shift(-2)  
df['t-1'] = df.t.shift(1)  
df['t-2'] = df.t.shift(2)
```

t-2	t-1	t	X	y
nan	nan	0	1	2
nan	0	1	2	3
0	1	2	3	4
1	2	3	4	5
2	3	4	5	6
3	4	5	6	7
4	5	6	7	8
5	6	7	8	9
6	7	8	9	nan
7	8	9	nan	nan

# TIME-SERIES TO SUPERVISED LEARNING

```
import pandas as pd
```

```
df = DataFrame()  
df['t'] = range(10)  
x, y = tseries_to_ml(  
    df, bsteps=2, fsteps=2)
```

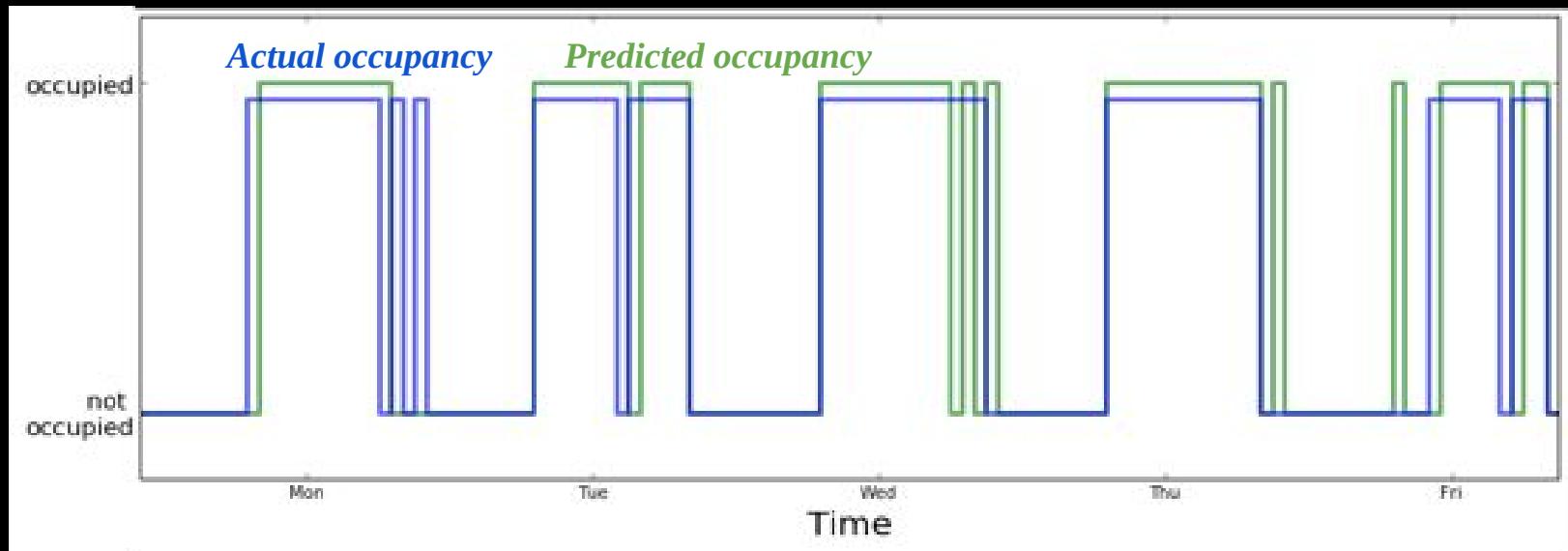
```
RandomForestClassifier().fit(x, y)
```

X

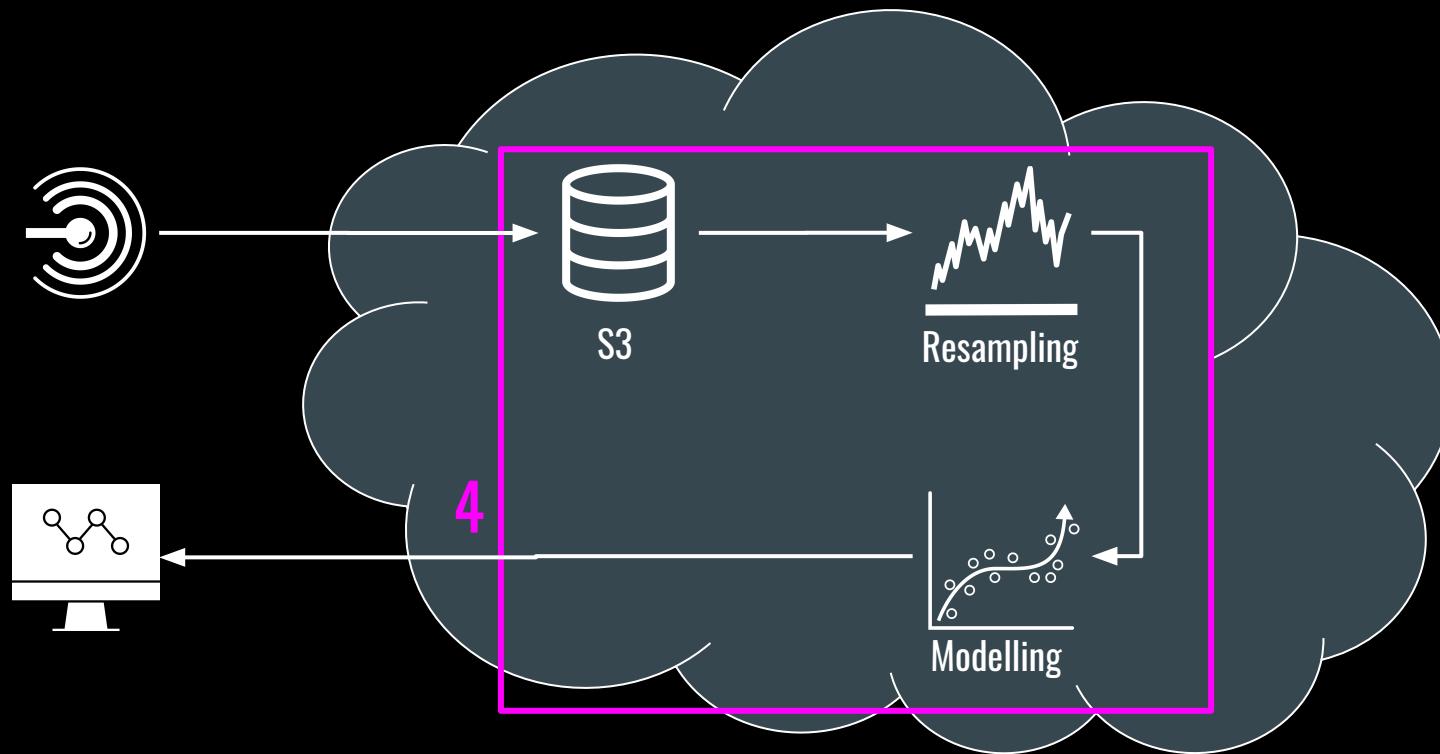
y

t-2	t-1	t	t+1	t+2
nan	nan	0	1	2
nan	0	1	2	3
0	1	2	3	4
1	2	3	4	5
2	3	4	5	6
3	4	5	6	7
4	5	6	7	8
5	6	7	8	9
6	7	8	9	nan
7	8	9	nan	nan

# OCCUPANCY PREDICTION

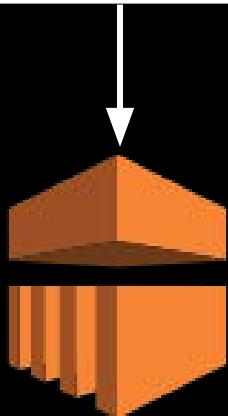


# 4 SCALING AND AUTOMATING



# SCALING AND AUTOMATING

FOR AUTOMATING



EMR



FOR SCALING



Spark

# SAMPLE SPARK APP



```
from spark import SparkContext
```

```
def main():
```

```
    sc = SparkContext(appName='TrainModels')
```

```
    room_sensors = grouped_sensors(gby='room', type='motion')
```

```
    rdd = sc.parallelize(room_sensors)
```

```
    rdd.map(load_sensor_data)
```

```
        .map(train_predictor)
```

```
    sc.stop()
```

```
if __name__ == '__main__':  
    main()
```

[ (room\_08, [ms\_12, ms\_14, m\_23]),  
 (room\_13, [ms\_15, ms\_18]), ... ]

# SAMPLE EMR JOB



```
import boto3

emr = boto3.client('emr')

emr.run_job_flow(
    Name='MyCluster',
    ReleaseLabel='emr-5.4.0',
    Instances={
        'MasterInstanceType': 'm3.xlarge',
        'SlaveInstanceType': 'm3.xlarge',
        'InstanceCount': 5,
    },
    BootstrapActions=BOOTSTRAP_ACTIONS,
    Steps=JOB_STEPS)
```

# JOB STEPS

```
JOB_STEPS = [ {  
    'Name': 'Resample Timeseries',  
    'HadoopJarStep': {  
        'Jar': 'command-runner.jar',  
        'Args': ['spark-submit',  
                '/home/hadoop/resample_tseries.py'],  
    },  
}, {  
    'Name': 'Train Occupancy Models',  
    'HadoopJarStep': {  
        'Jar': 'command-runner.jar',  
        'Args': ['spark-submit',  
                '/home/hadoop/train_occupancy.py'],  
    },  
}]
```

Resamples data per sensor



Trains predictor per room



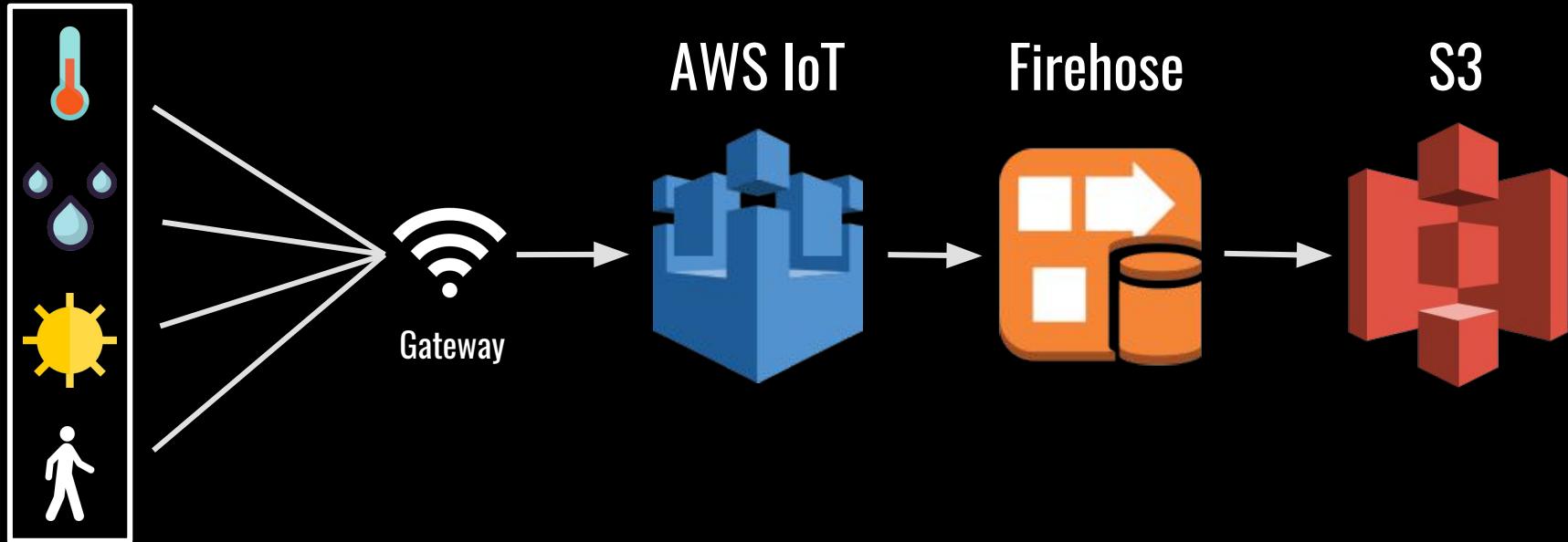
# BOOTSTRAP ACTIONS

```
BOOTSTRAP_ACTIONS = [{  
    'Name': 'Install Python Libraries',  
    'ScriptBootstrapAction': {  
        'Path': 's3://path/to/code/install_libraries.sh',  
        'Args': []  
    }  
}]
```

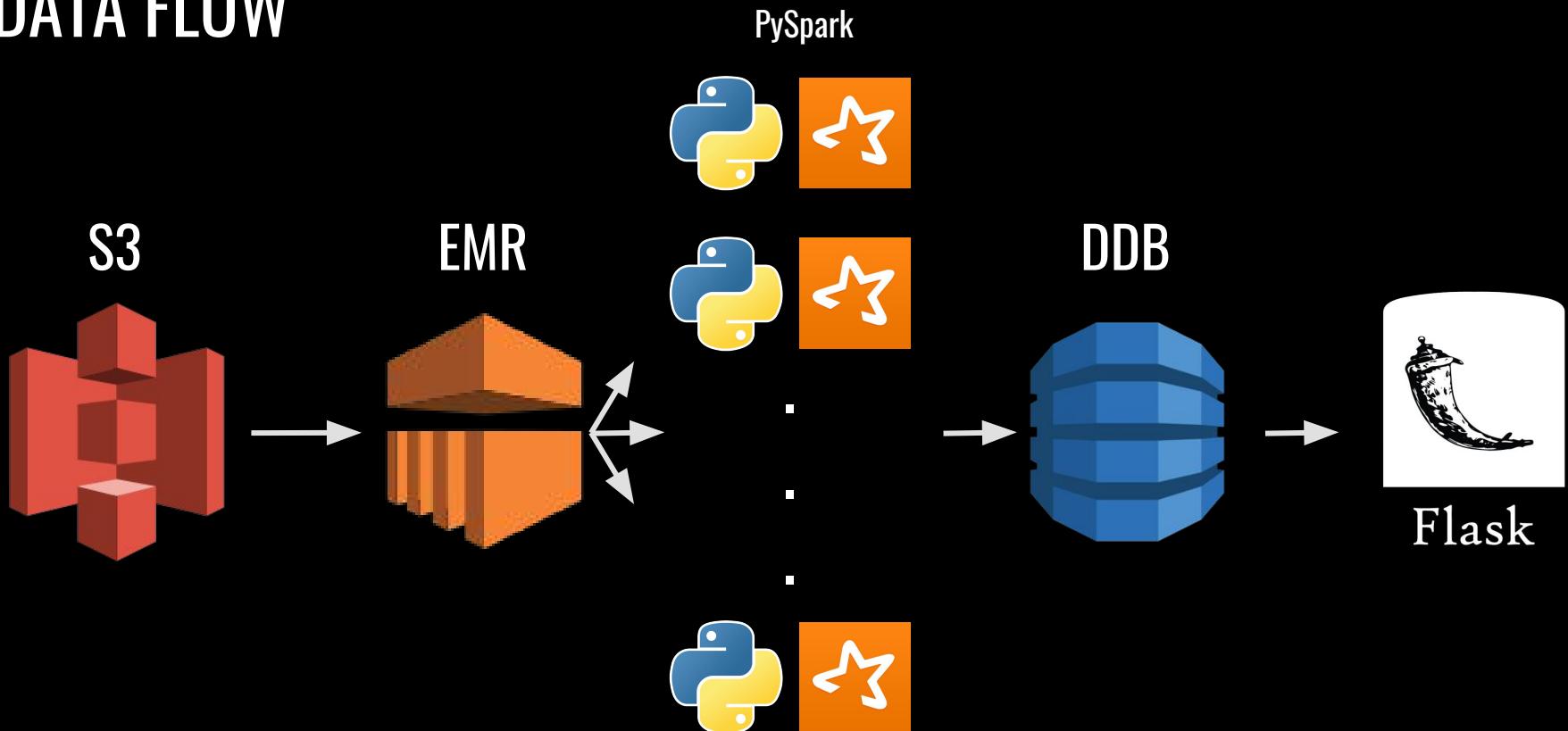


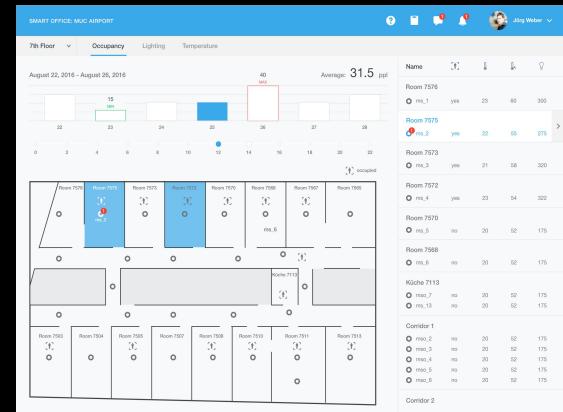
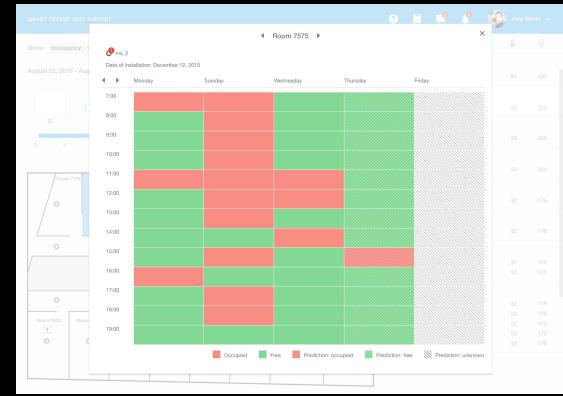
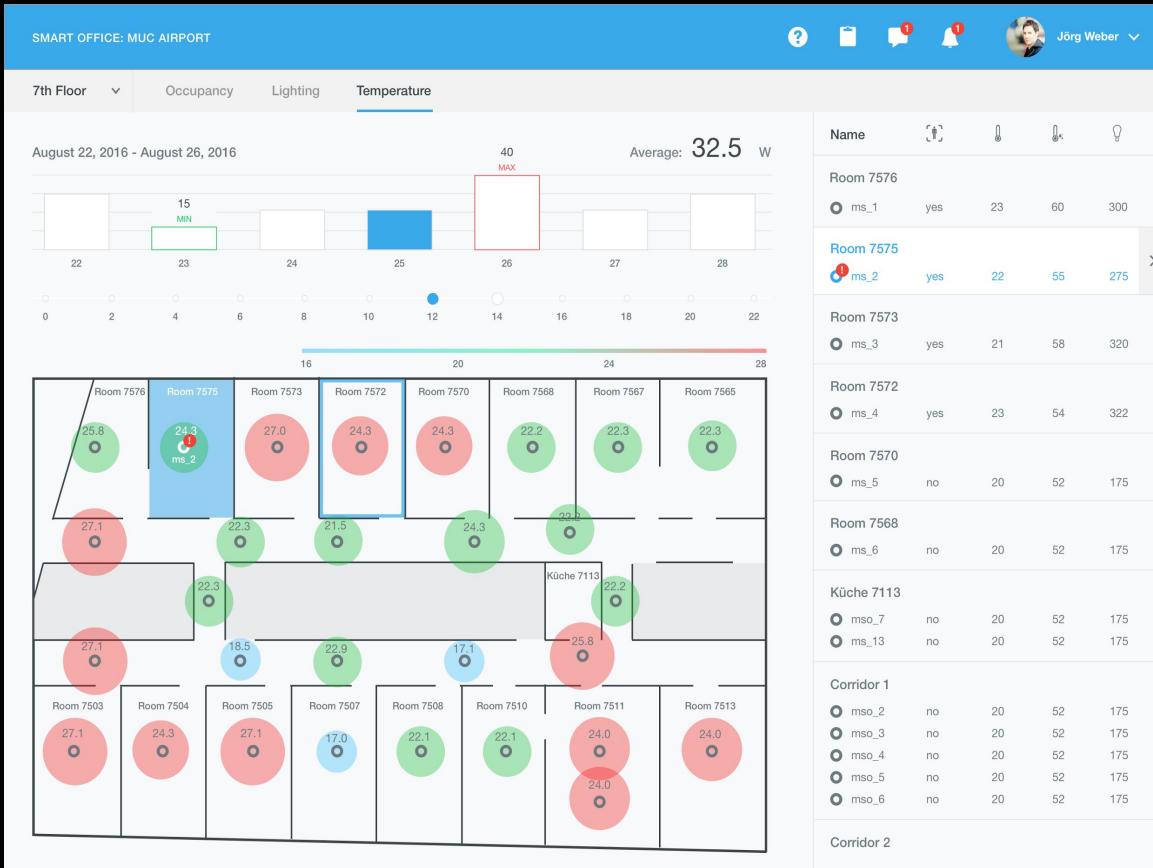
pip install numpy pandas scikit-learn

# DATA FLOW



# DATA FLOW







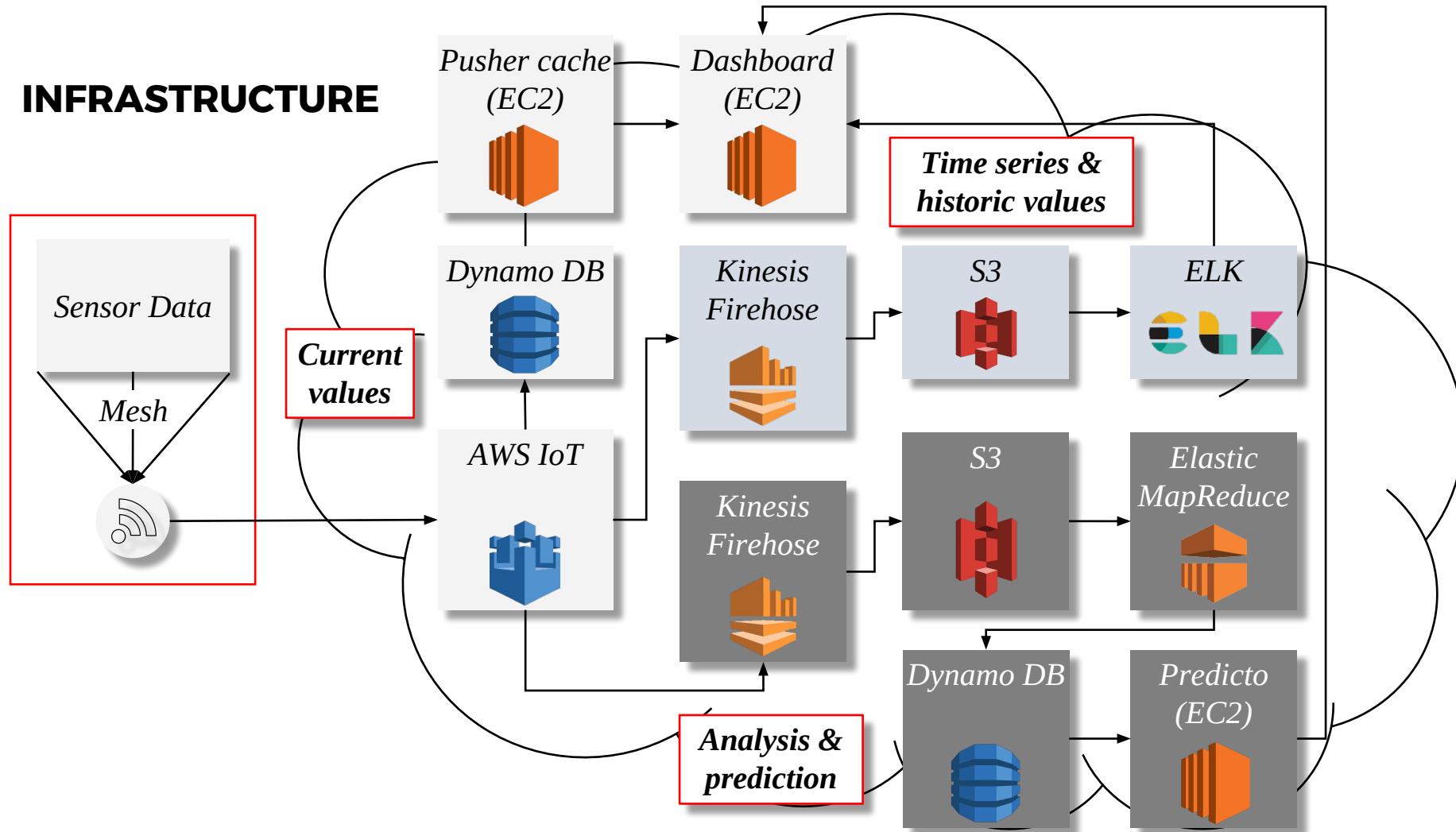
[www.wattx.io](http://www.wattx.io)

Rafael Schultze-Kraft | Data Scientist | rafael@wattx.io

# OUTLOOK

- ❑ Going Real-time: Stream Processing
- ❑ Going Big: SparkML vs Scikit-Learn
- ❑ Modelling: MCMC, LSTMs, ...
- ❑ Closing the Loop: Control of an Actuator

# INFRASTRUCTURE



# OCCUPANCY PREDICTION

