be fast, be agile, At Scale

# Appagile Data Science Workstation - Analytics from the cloud
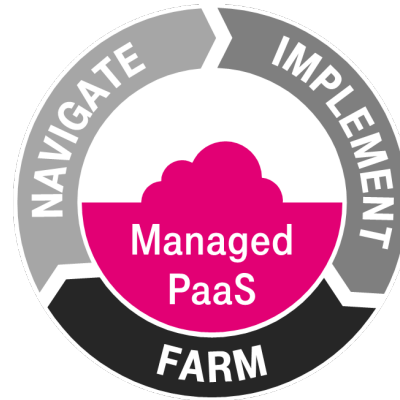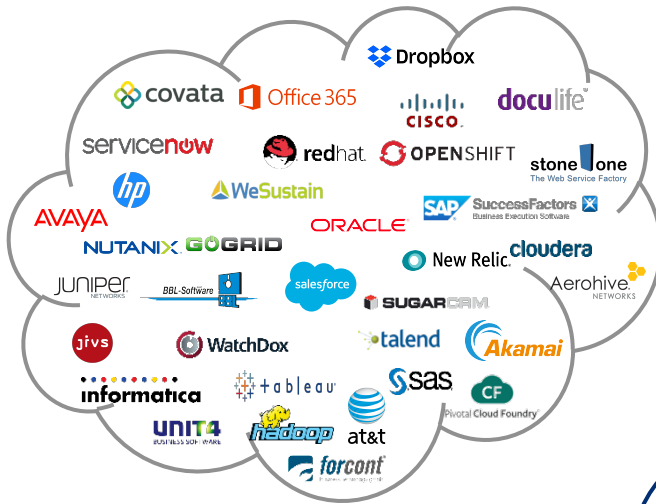
**Stefan Zosel**

May, 2018 – Openshift Anwenderforum
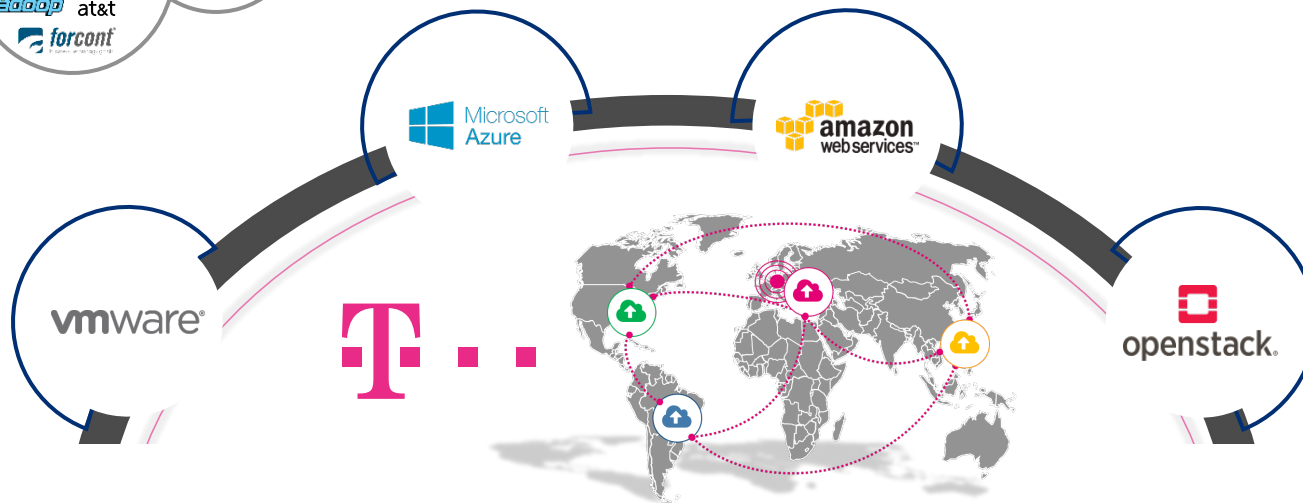
**T··Systems·**

# Digitization with T-Systems

Managed Cloud Services

Security approved operation

Openshift for Enterprises

# Warum DataScience auf Openshift?

Innovationen aus der Cloud

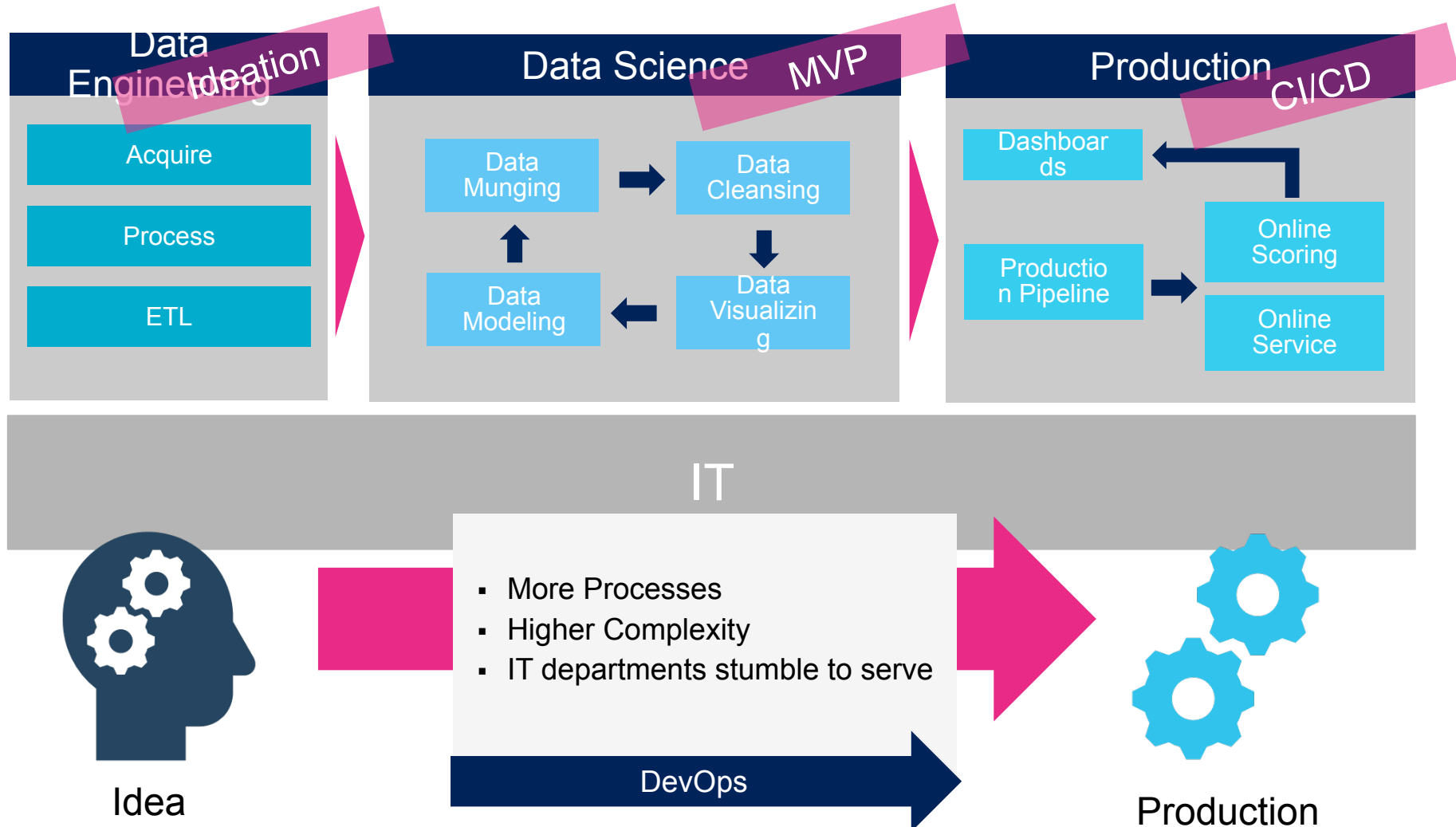Neue UseCases

AI/ML

BIG Data trifft PaaS

Neue Methoden

BIG Data trifft Cloud

Viele neue SW & Services

Self Service Bedarf

# Data Science - DevOps
## From data to product

**Data Engineering** — *Ideation*

- Acquire
- Process
- ETL

**Data Science** — *MVP*

Data Munging → Data Cleansing → Data Visualizing → Data Modeling → Data Munging

**Production** — *CI/CD*

- Dashboards
- Production Pipeline → Online Scoring
- Online Service

**IT**

Idea

- More Processes
- Higher Complexity
- IT departments stumble to serve

DevOps

Production

**T··Systems·**

# Data Science DevOps – not Admin!
## The pain points

**Sound Similar to Developer**

**Setup**

- Provision of environments takes forever
- Compatibility problems

- Provision of data access takes time
- Connection to data sources is cumbersome

**Connection**

**Scalability**

- Hardware resources on desktops limited
- Models can only be trained and tested on samples
- Popular frameworks are not parallelized.

- Language and version conflicts
- Sharing results takes extra effort

**Sharing**

# Data Science DEVOPS
## What else Would be nice?

**Dockerization in the cloud**

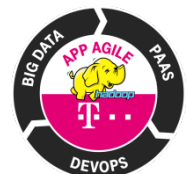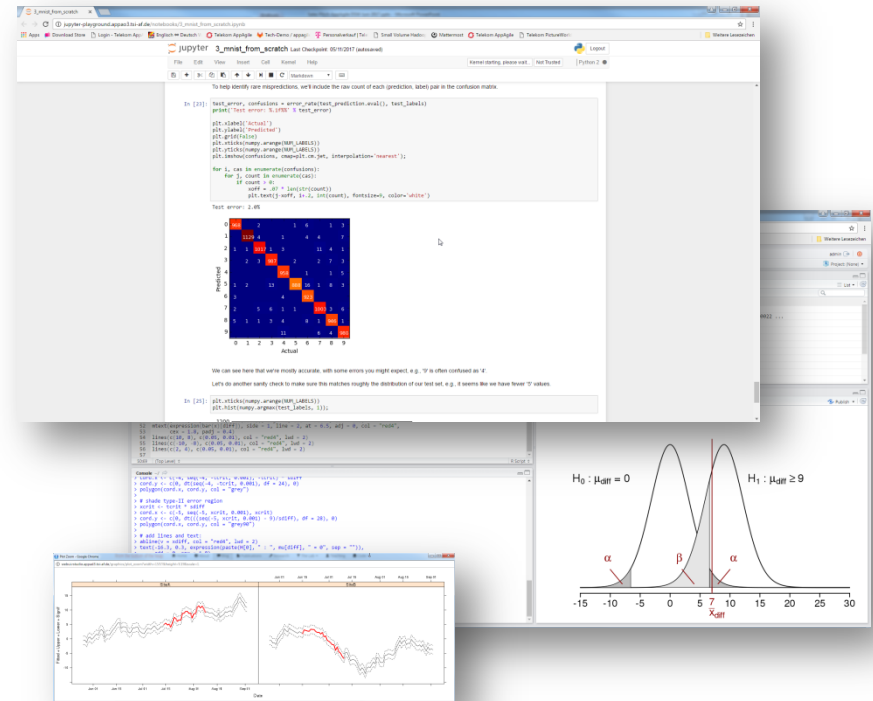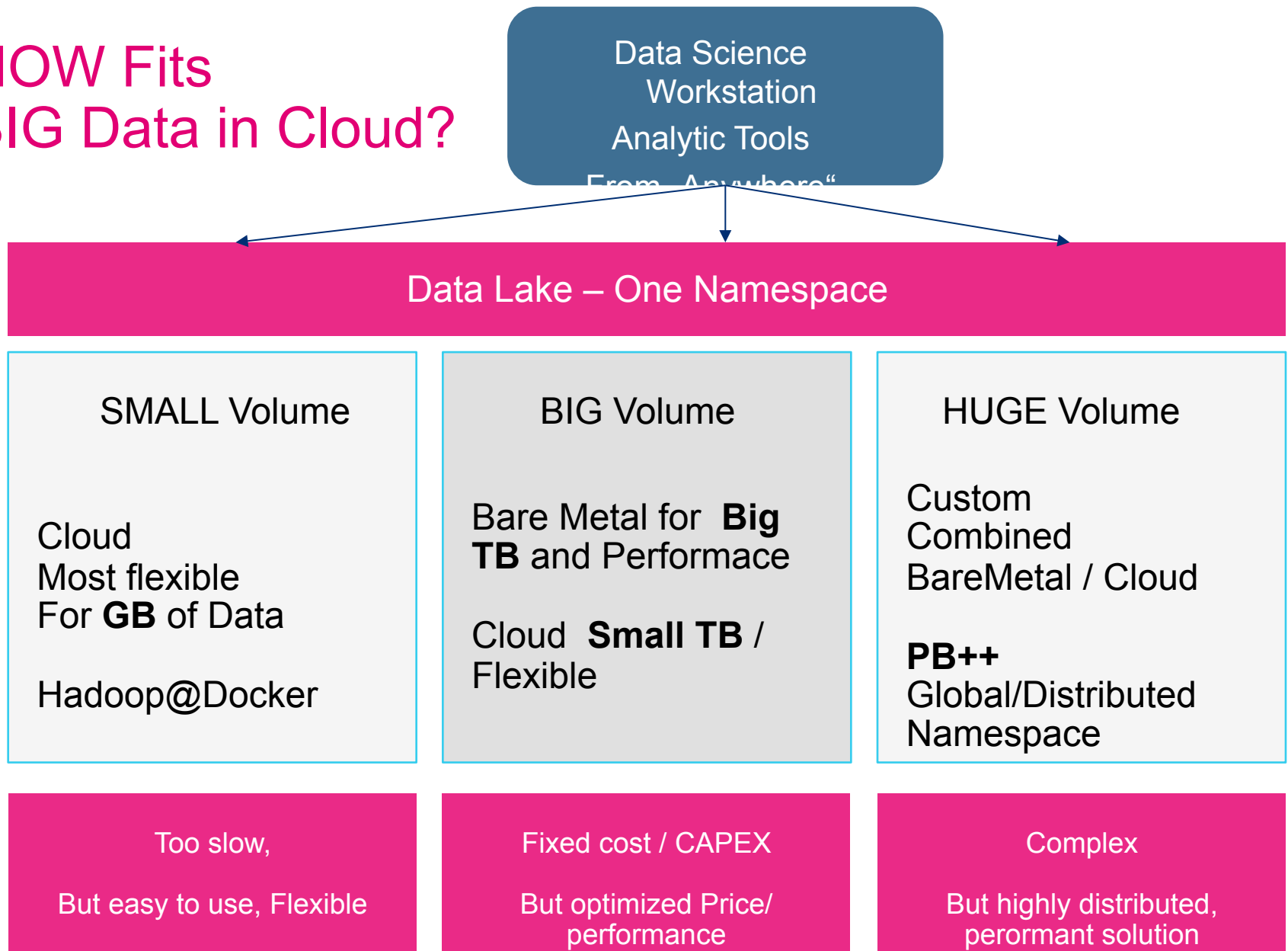| Sandboxing | PoC | Training |
|---|---|---|
| **Try** new ideas **without destroying** anything | Develop your **PoC** agile, **at fast pace** in the cloud | **Train** your team in the cloud on **new ML frameworks** |
| Locality | Ad-hoc Analysis | CI/CD |
| Run **analyses** in the cloud where your **data is located** | **Set up quickly** an environment **with all tools** | Work in **parallel** and **contiuously** on project **modules** |

# Presenting Appagile's Data science Workstation
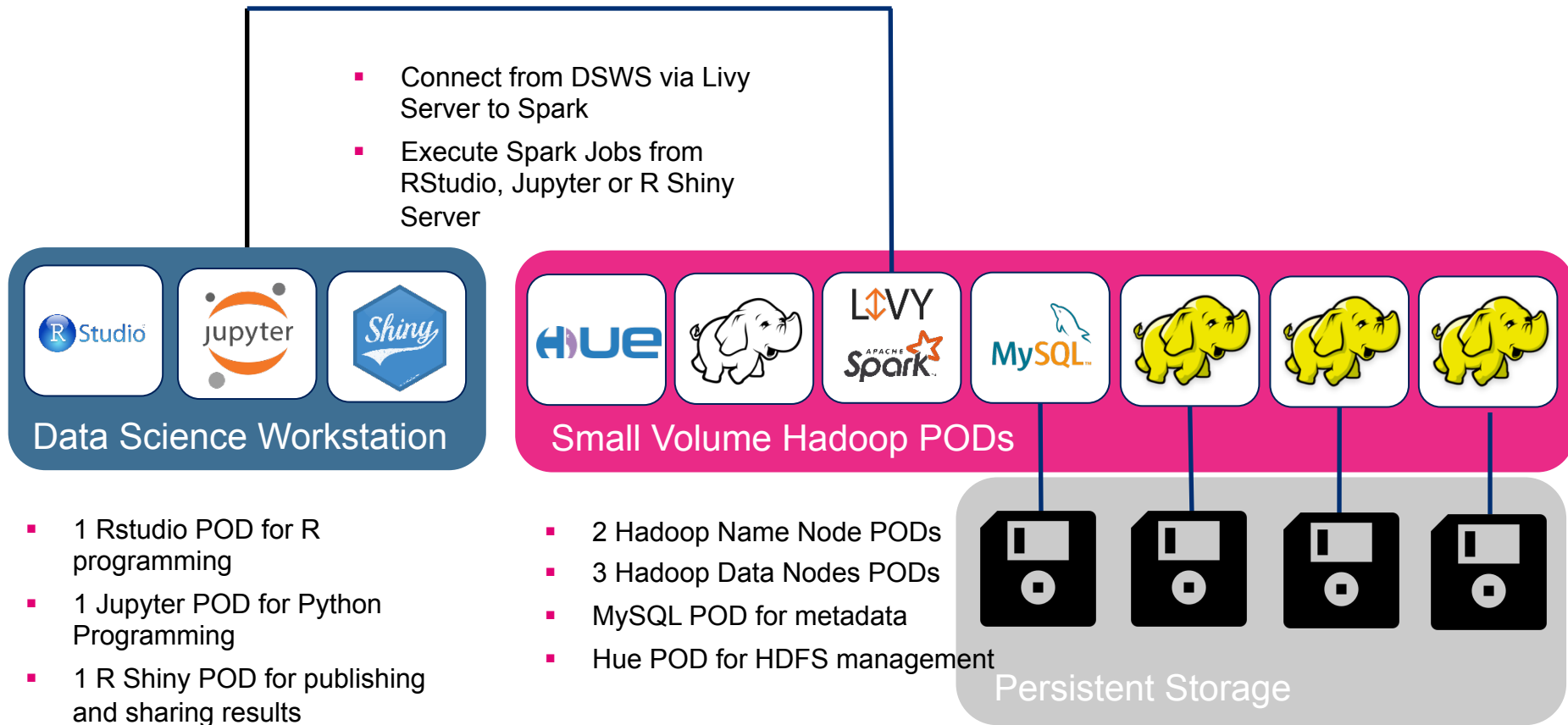## be fast, be agile, at scale

- Quick and easy self-service
- **Dedicated** Data Science Engineer
- Working with **AppAgile's Small Volume Hadoop** cluster
- R, Python and Scala in multiple versions
- Common machine learning **frameworks**
- **Visualization** via notebooks, markdown and R Shiny Server
- Available on **OTC**, **vCloud**, and **Microsoft Azure**
- From sample to big data

# HOW Fits
# BIG Data in Cloud?

Data Science
Workstation

Analytic Tools

From "Anywhere"

## Data Lake – One Namespace

| SMALL Volume | BIG Volume | HUGE Volume |
|---|---|---|
| Cloud<br>Most flexible<br>For **GB** of Data<br><br>Hadoop@Docker | Bare Metal for **Big TB** and Performace<br><br>Cloud **Small TB** / Flexible | Custom Combined BareMetal / Cloud<br><br>**PB++**<br>Global/Distributed Namespace |
| Too slow,<br><br>But easy to use, Flexible | Fixed cost / CAPEX<br><br>But optimized Price/ performance | Complex<br><br>But highly distributed, perormant solution |

**T··Systems·**

# Appagile Small volume hadoop
## data Science Setup Example

- Connect from DSWS via Livy Server to Spark
- Execute Spark Jobs from RStudio, Jupyter or R Shiny Server



**Data Science Workstation**

**Small Volume Hadoop PODs**

**Persistent Storage**

- 1 Rstudio POD for R programming
- 1 Jupyter POD for Python Programming
- 1 R Shiny POD for publishing and sharing results

- 2 Hadoop Name Node PODs
- 3 Hadoop Data Nodes PODs
- MySQL POD for metadata
- Hue POD for HDFS management

**Infrastructure: Private vCloud / Public vCloud/ OTC / Microsoft Azure**

# Some Screens

Projects

Overview

Applications

Builds

Resources

Storage

Monitoring

∨ APPAGILE HUE

https://hue.hado

hadoop-hue

Deployment appagile-hue – 7 days ago                                    #2

CONTAINER: APPAGILE-HUE

**Image:** smalldata/appagile-hue
**Ports:** 8888/TCP

1
pod

No grouped services.

No services are grouped with hadoop-hue. Add a service to group the

**Group Service**

∨ HADOOP QUICKSTART

hadoop-local

Deployment appagile-hadoop-master – 7 days ago                          #1

CONTAINER: APPAGILE-HADOOP-MASTER

**Image:** smalldata/appagile-hadoop
**Ports:** 22/TCP , 6066/TCP , 7077/TCP ,
8030/TCP , 8031/TCP , 8032/TCP , 8033/TCP ,
8080/TCP , 8088/TCP , 8998/TCP , 9000/TCP ,
9083/TCP , 10000/TCP , 10002/TCP ,
19888/TCP , 50070/TCP , 50090/TCP

1
pod

No grouped services.

No services are grouped with hadoop-local. Add a service to group the

**Group Service**

hadoop-local

Deployment appagile-hadoop-slave1 – 7 days ago                          #1

CONTAINER: APPAGILE-HADOOP-SLAVE1

**Image:** smalldata/appagile-hadoop
**Ports:** 22/TCP , 8040/TCP , 8042/TCP ,
8081/TCP , 50010/TCP , 50020/TCP ,
50075/TCP

1
pod

```
In [22]:  K = 2
          nearest_partition = np.argpartition(dist_sq, K + 1, axis=1)
```

In order to visualize this network of neighbors, let's quickly plot the points along with lines representing the connections from each point to its two nearest neighbors:

```
In [23]:  plt.scatter(X[:, 0], X[:, 1], s=100)

          # draw lines from each point to its two nearest neighbors
          K = 2

          for i in range(X.shape[0]):
              for j in nearest_partition[i, :K+1]:
                  # plot a line from X[i] to X[j]
                  # use some zip magic to make it happen:
                  plt.plot(*zip(X[j], X[i]), color='black')
```



Each point in the plot has lines drawn to its two nearest neighbors. At first glance, it might seem strange that some of the points have more than two lines coming out of them: this is due to the fact that if point A is one of the two nearest neighbors of point B, this does not necessarily imply that point B is one of the two nearest neighbors of point A.

Although the broadcasting and row-wise sorting of this approach might seem less straightforward than writing a loop, it turns out to be a very efficient way of operating on this data in Python. You might be tempted to do the same type of operation by manually looping through the data and sorting each set of neighbors individually, but this would almost certainly lead to a slower algorithm than the vectorized version we used. The beauty of this approach is that it's written in a way that's agnostic to the size of the input data: we could just as easily compute the neighbors among 100 or 1,000,000 points in any number of dimensions, and the code would look the same.

Finally, I'll note that when doing very large nearest neighbor searches, there are tree-based and/or approximate algorithms that can scale as $\mathcal{O}[N \log N]$ or better rather than the $\mathcal{O}[N^2]$ of the brute-force algorithm. One example of this is the KD-Tree, implemented in Scikit-learn.

## Aside: Big-O Notation

Search ...

**GROUPS**

Newest | Active | Popular

PaaS Q&A
active 34 minutes ago

Data Science Work-bench
active 11 hours, 8 minutes ago

BIG Data
active 5 months, 4 weeks ago

Imprint

# DATA SCIENCE WORKBENCH



Leave Group

Private Group   11 hours ago
This community supports data science to work with our analytic workbench, delivered per data science workbench.

**Group Admins**

| Home | RSS | Show: — Everything — |
|------|-----|---------------------|

# Community - AppAgile.IO

# Appagile data science workstation
## What's in the Package:

**Data Science Engineer**

- For named users
- Dedicated to DS support
- Building of DS libraries
- In charge of DS community
- DS Q&A (Best Practise)

**Analytics Services – Data Science Workstation**

- Ready-to-use workstation for DS
- DS self-service
- Predefined data and analytic modules

Dedicated DS Engineer

Container Images

**Data Science Community**

Middleware - Small Volume Hadoop

- Hadoop in dockers on OpenShift
- Horizontally scalable on demand
- Free assignment of capacities per node
- On public and private cloud
- Compatible to big data solution

**Container Repository**

**Data Science Workstation:**
- For named user

**Small Volume:**
- Unmanaged - Free