



Distributed DataFrame (DDF)

Simplifying Big Data
For The Rest Of Us

Christopher Nguyen, PhD
Co-Founder & CEO

Agenda

- I. Motivation: Problem & Solution**
- 2. DDF Features & Benefits**
- 3. Demo**



Christopher Nguyen, PhD
Adatao Inc.
Co-Founder & CEO

- Former **Engineering Director of Google Apps** (Google Founders' Award)
- Former Professor and Co-Founder of the **Computer Engineering program at HKUST**
- **PhD Stanford, BS U.C. Berkeley** *Summa cum Laude*

A vibrant photograph of seven diverse children of various ethnicities running together in a lush green field dotted with bright yellow flowers. The children are dressed in casual summer clothing like t-shirts, shorts, and sunglasses. They are all smiling and appear to be having a great time. The background features a clear blue sky and a line of trees, suggesting a rural or park setting.

**small data
world**

BIG BAD DATA

Pig
HDFS
HBASE
Hive
Cascading
MapReduce
SummingBird
Crunch
Storm
Scalding
RHadoop

BIG BAD DATA



```
class RandomSplitRDDC(prev: RDD[Array[Double]], seed: Long,  
upper: Double, isTraining: Boolean)  
extends RDD[Array[Double]](prev) {  
override def getPartitions: Array[Partition] = {  
  val rg = new Random(seed)  
  firstParent[Array[Double]].partitions.map(x => new  
    Partition(x.id, rg.nextInt))  
}  
override def getPreferredLocations(split: Partition):  
  Iterator[String] = firstParent[Array[Double]].preferredLocations(split.asIn-  
  n].prev  
override def compute(splitIn: Partition, context: TaskContext):  
  Iterator[Array[Double]] = {  
  val split = splitIn.asInstanceOf[SeededPartition]  
  val rand = new Random(split.seed)  
  if (isTraining) {  
    firstParent[Array[Double]].iterator(split.prev, context).filter(x =>  
      val z = rand.nextDouble;  
      z < lower || z >= upper  
    )  
  } else {  
    firstParent[Array[Double]].iterator(split.prev, context).filter(x =>  
      val z = rand.nextDouble;  
      lower <= z && z < upper  
    )  
  }  
}  
def randomSplit[Array[Double]](rdd: RDD[Array[Double]], numSplits: Int,  
  trainingSize: Double, seed: Long): Iterator[(RDD[Array[Double]],  
  RDD[Array[Double]])] = {  
  require(0 < trainingSize && trainingSize < 1)  
  val rg = new Random(seed)  
  (1 to numSplits).map(_ => rg.nextInt).map(z =>  
    (new RandomSplitRDD(rdd, z, 0, 1.0 - trainingSize, true),  
     new RandomSplitRDD(rdd, z, 0, 1.0 - trainingSize, false))).toIterator  
}
```



**WHY
Can't It Be
Even
Simpler?**



What Happiness Looks Like

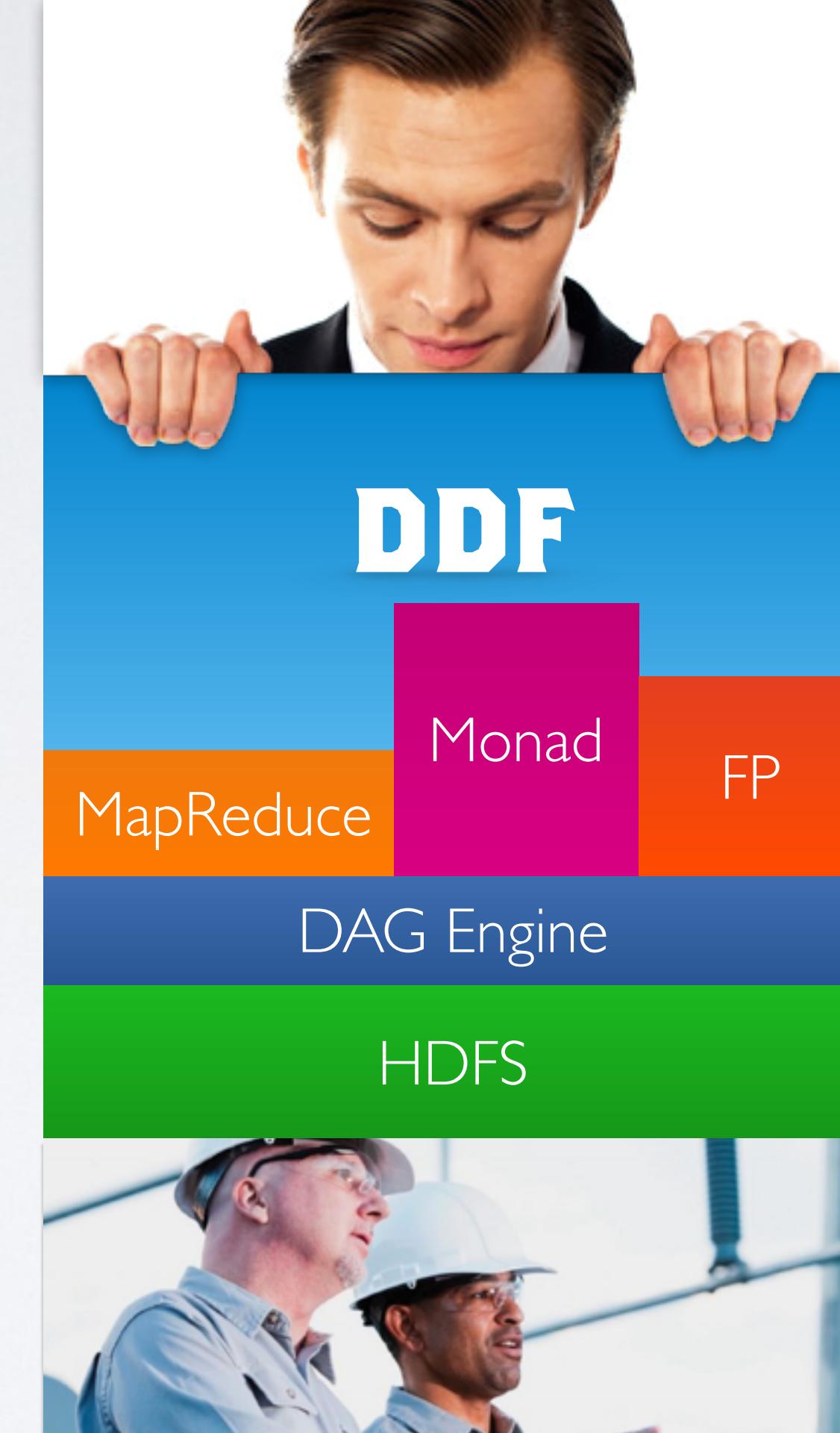
```
class RandomSplitRDD( prev: RDD[Array[Double]], seed: Long  
upper: Double, isTraining: Boolean)  
extends RDD[Array[Double]](prev) {  
override def getPartitions: Array[Partition] = {  
  val rg = new Random(seed)  
  firstParent[Array[Double]].partitions.map(x => new  
    rg.nextInt))  
}  
override def getPreferredLocations(split: Partition):  
  firstParent[Array[Double]].preferredLocations(split.asIn  
n].prev)  
override def compute(splitIn: Partition, context: TaskC  
Iterator[Array[Double]] = {  
  val split = splitIn.asInstanceOf[SeededPartition]  
  val rand = new Random(split.seed)  
  if (isTraining) {  
    firstParent[Array[Double]].iterator(split.prev, context).fir  
    val z = rand.nextDouble;  
    z < lower || z >= upper  
  }  
  } else {  
    firstParent[Array[Double]].iterator(split.prev, context).filter(x => {  
      val z = rand.nextDouble;  
      lower <= z && z < upper  
    })  
  }  
}  
def randomSplit[Array[Double]](rdd: RDD[Array[Double]], numSplits: Int,  
  trainingSize: Double, seed: Long): Iterator[(RDD[Array[Double]],  
  RDD[Array[Double]])] = {  
  require(0 < trainingSize && trainingSize < 1)  
  val rg = new Random(seed)  
  (1 to numSplits).map(_ => rg.nextInt).map(z =>  
    (new RandomSplitRDD(rdd, z, 0, 1.0 - trainingSize, true),  
     new RandomSplitRDD(rdd, z, 0, 1.0 - trainingSize, false))).toIterator  
}
```

```
table = load("housePrices")  
table.dropNA  
table.train.glm("price", "sf,zip,beds,baths")
```

Top-Down

vs.

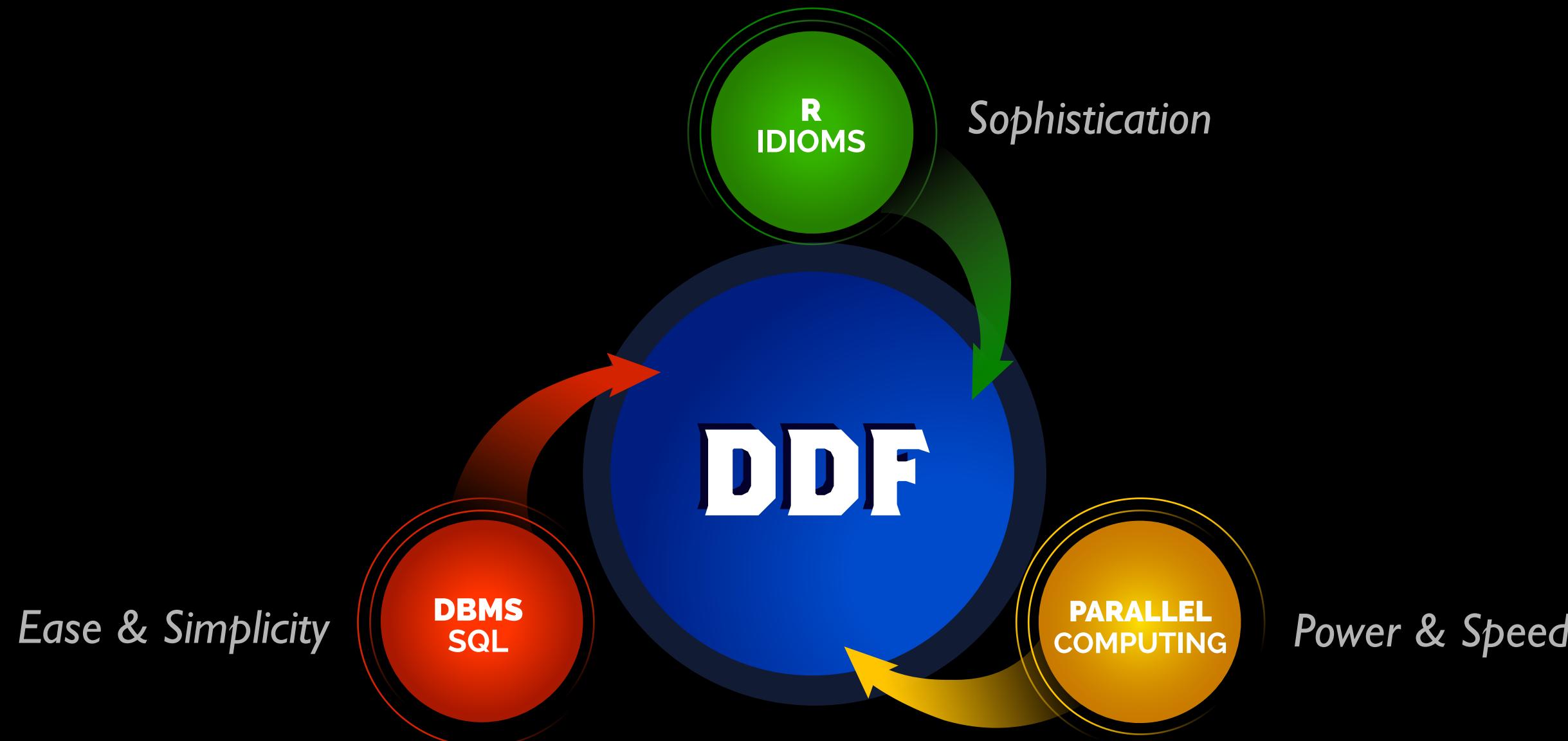
Bottom-Up



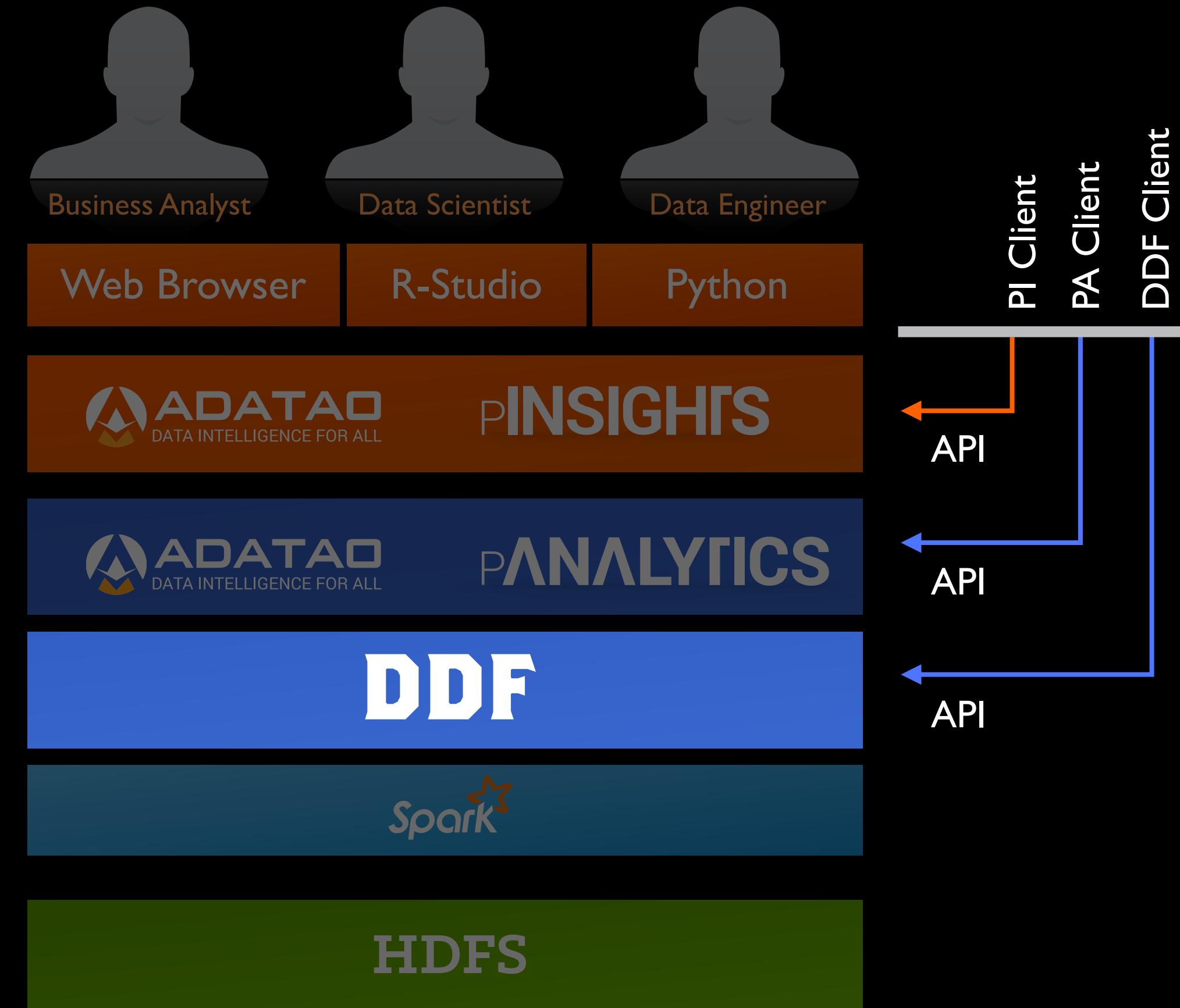
TOP SECRET



Design Principles



Example:



Demo Deployment Diagram



DDF Offers



1

Table-Like Abstraction on Top of Big Data



6

DDF Offers



2

Native R Data.frame Experience

```
ddf.getSummary()  
ddf.getFiveNum()  
ddf.binning("distance", 10)  
ddf.scale(SCALE.MINMAX)
```

DDF Offers



3

Focus on Analytics, Not MapReduce

Data
Wrangling

Model
Building & Validation

Real-Time
Scoring

`ddf.dropNA()`

`ddf.transform("speed
=distance/duration")`

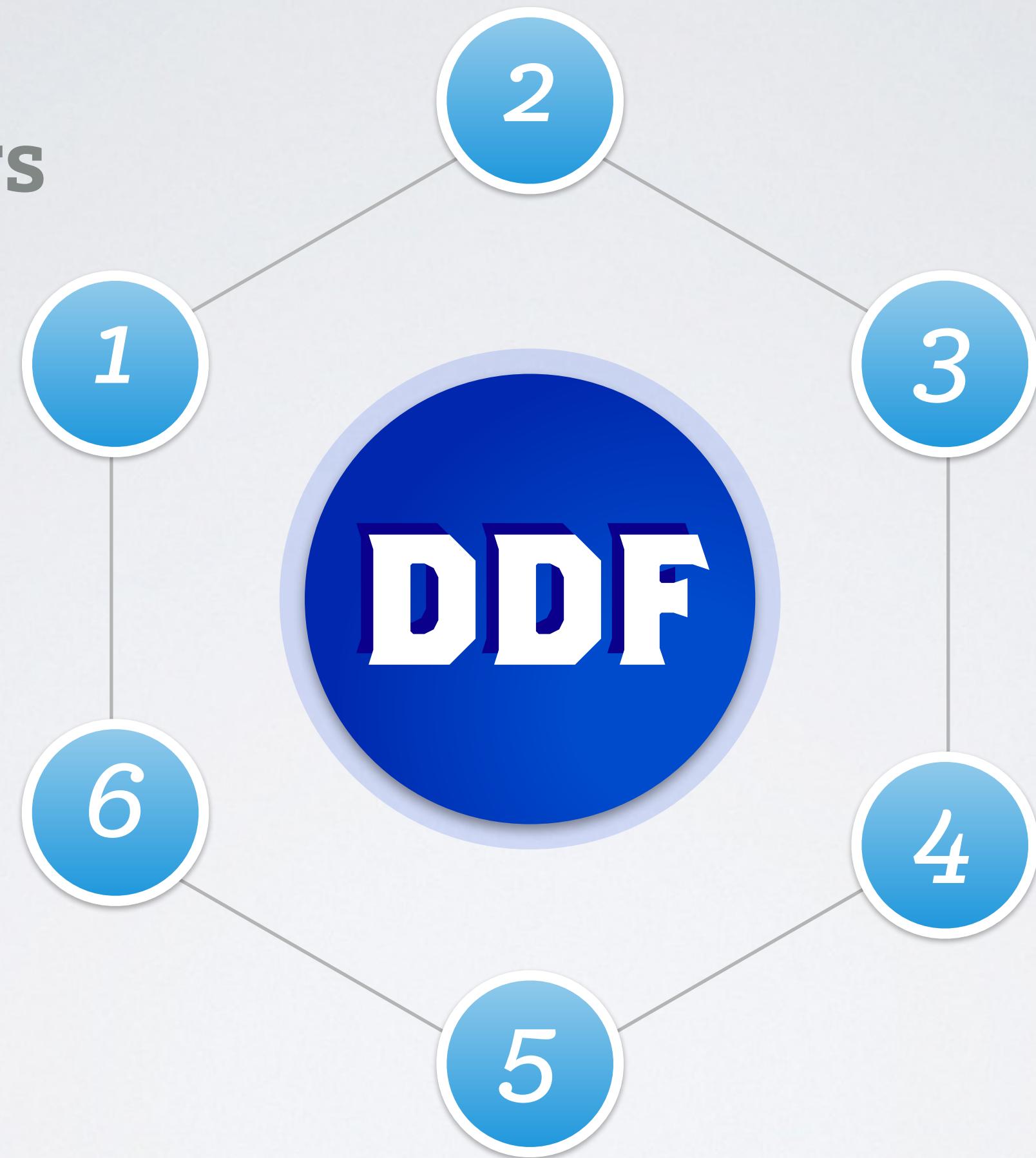
`ddf.lm(0.1, 10)`

`ddf.roc(testddf)`

`manager.loadModel(uri)`

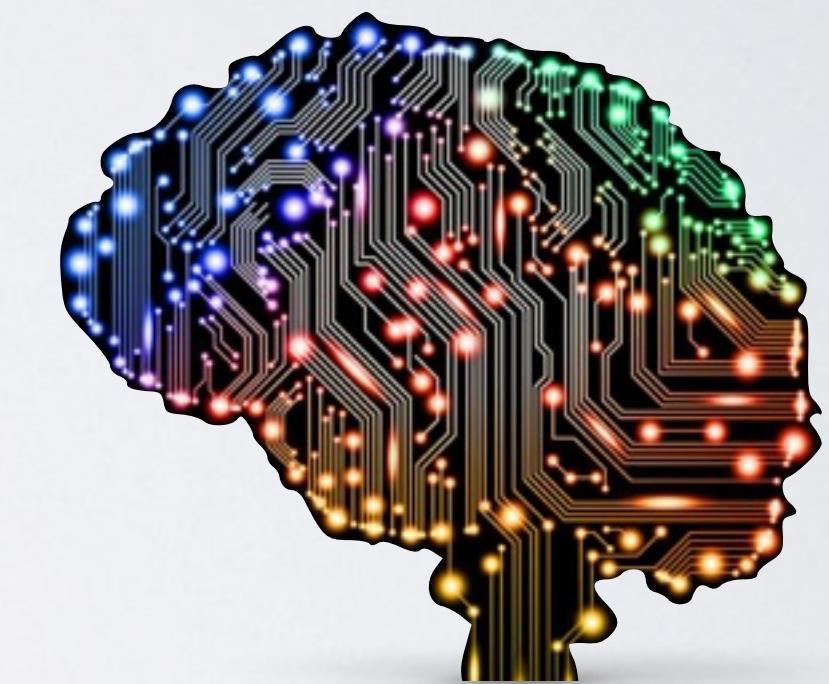
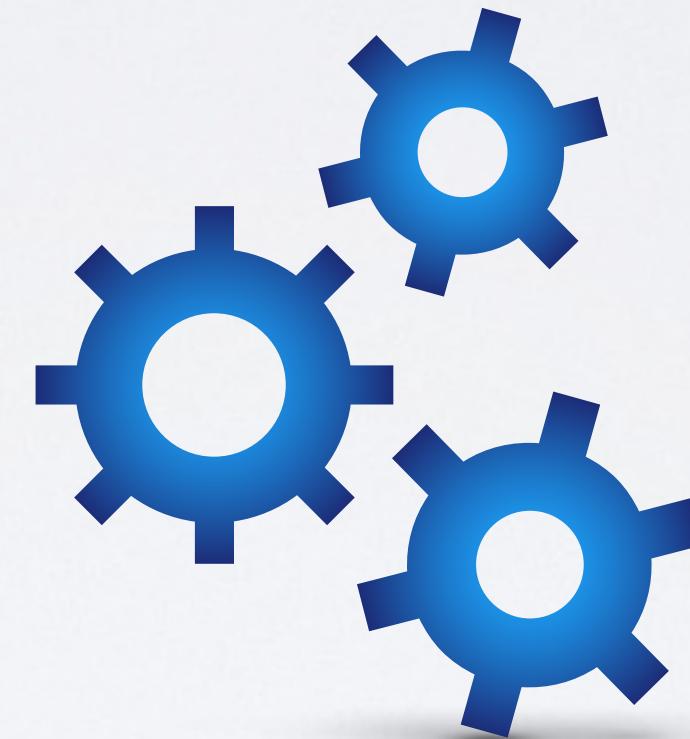
`lm.predict(point)`

DDF Offers



4

Seamlessly Integrate with External ML Libraries



Data

Algorithm

Intelligence

DDF Offers



5

Share DDFs Seamlessly & Efficiently



*Can I see
your data?*



<ddf://adatao/airline>

DDF Offers





Work in Your Preferred Language



python



...even Plain English!



Table-Like Abstraction
on Top of Big Data

**Work in Your
Preferred Language**

Share DDFs
Seamlessly & efficiently





DDF is to
Big Data
as
SQL is to
Small Data



We're Open Sourcing DDF!!!!

1. Accept Core Committers
2. Clean up & Prep
3. Open up for public access

wwwADATAO.com/ddf

