



What's New for Machine Learning with Oracle Database and Hadoop

Eric Grancher

Manuel Martín Márquez

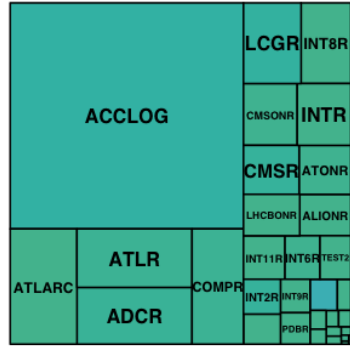


Outline

- CERN Data environment
- Introduction to Machine Learning
- Our Machine Learning path
 - Introduction to a real use case
 - R – prototyping and validating ideas
 - First Scalable attend – Oracle R Enterprise
 - Hadoop and the analytic transformation
 - Oracle Advance Analytics for Hadoop
 - Spark
 - TensorFlow + Spark

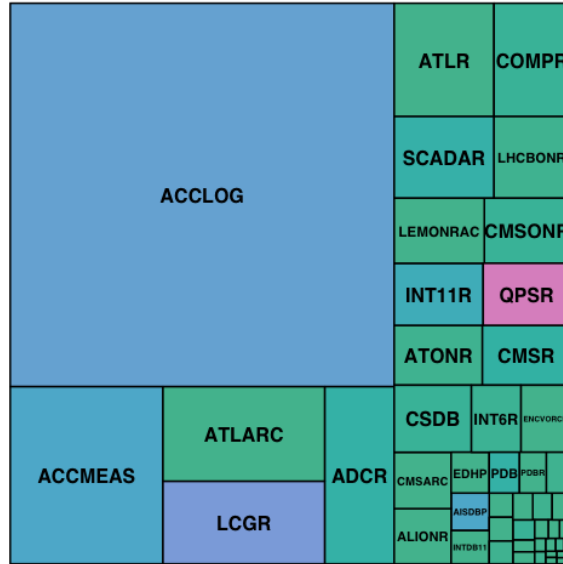
CERN Database Environment

CERN databases 201210,
area is size (total=299TB)



0 5 10 15 20 25
TB redo per month, sum=120 TB

CERN databases 201512,
area is size (total=750TB), color is redo activity

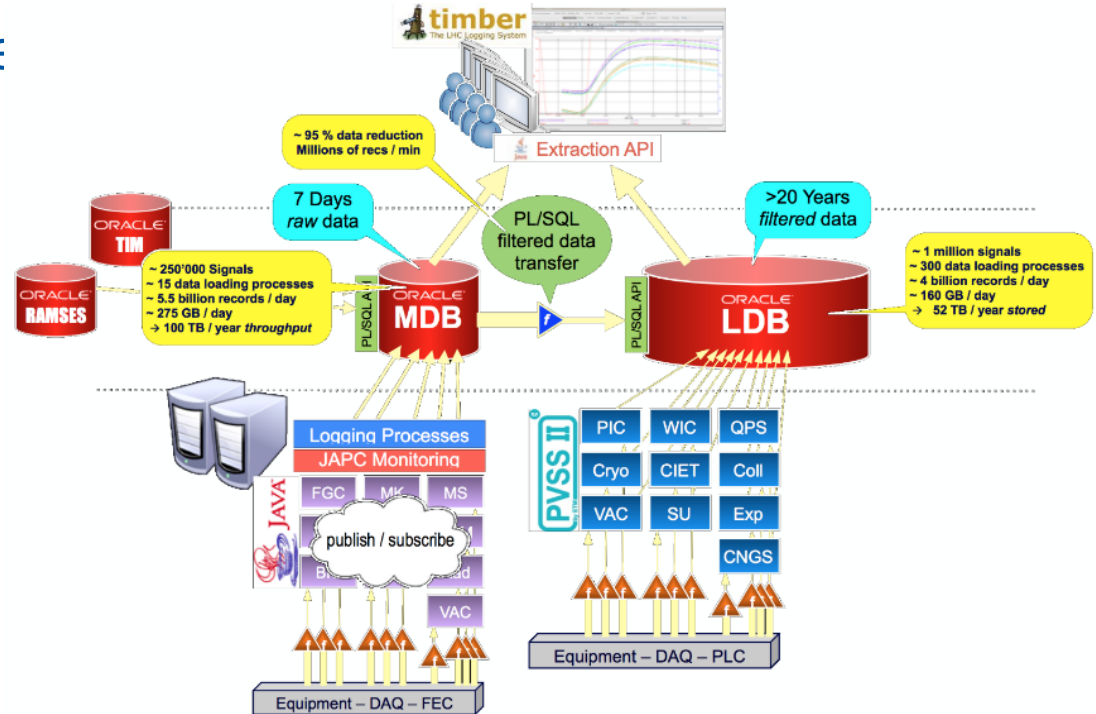


0 10 20 30 40 50 60 70 80 90 100 110
TB redo per month, sum=494 TB

	October 2012	December 2015
Max size	ACCLOG 136TB	ACCLOG 352TB
Max redo	ACCMEAS 27TB / month	QPSR 115TB / month

CERN Control Systems

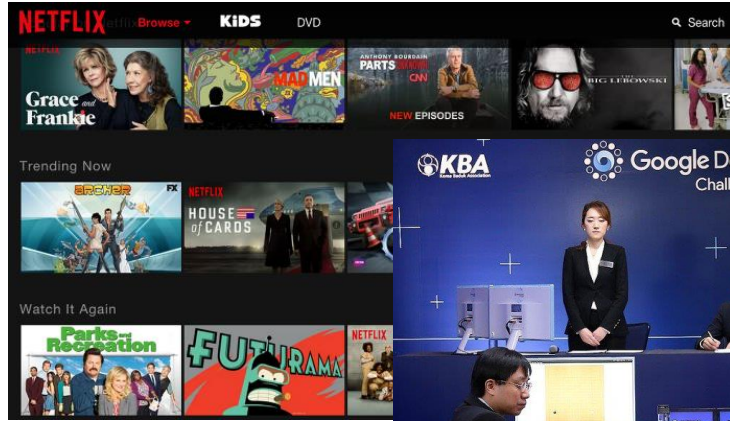
- IoT and Control Systems
 - Cryogenics
 - Vacuum
 - Machine Protection
 - Power Converters
 - QPS
- Accelerator Logging Service
 - ~ 275 GB/day
 - Storing more than 50 TB / year
 - Data acquisition
 - CERN accelerator complex
 - Related subsystems
 - Experiments
 - Around 1 million signals



Machine Learning (ML)

- ML is a branch of artificial intelligence:
 - Uses computing based systems to make sense out of data
 - Extracting patterns, fitting data to functions, classifying data, etc
 - ML systems can learn and improve
 - With historical data, time and experience
 - Bridges theoretical computer science and real noise data.

ML in real-life



Higgs challenge

Completed • \$13,000 • 1,785 teams

Higgs Boson Machine Learning Challenge

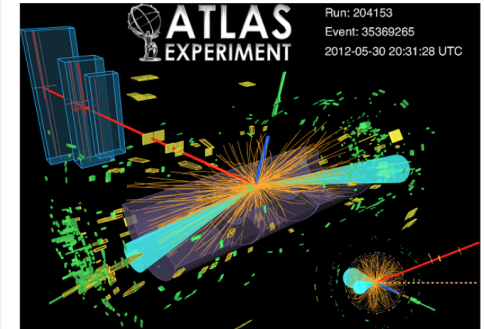
Mon 12 May 2014 - Mon 15 Sep 2014 (2 years ago)

Dashboard

- Home
- Data
- Make a submission
- Information
 - Description
 - Evaluation
 - Rules
 - Prizes
 - About the Sponsors
 - Timeline
 - Winners
- Forum
- Leaderboard
 - Public
 - Private
- My Team
 - Your model
 - GitHub
- My Submissions

Competition Details » [Get the Data](#) » [Make a submission](#)

Use the ATLAS experiment to identify the Higgs boson

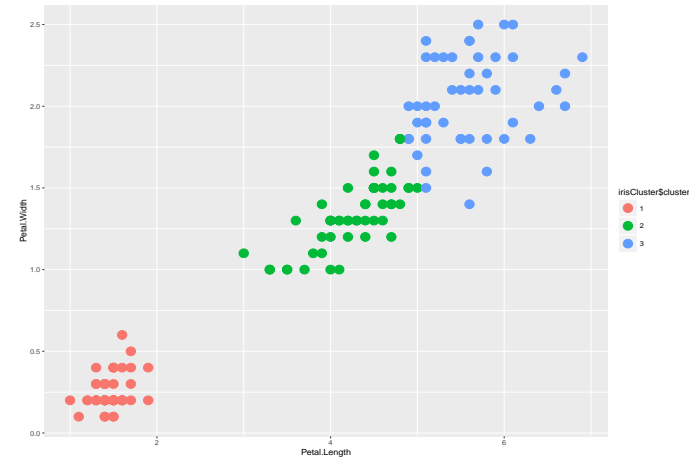


Discovery of the long awaited Higgs boson was announced July 4, 2012 and confirmed

Supervised and Unsupervised Learning

- Unsupervised Learning
 - There are not predefined and known set of outcomes
 - Look for hidden patterns and relations in the data
 - A typical example: Clustering

	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width
1	5.1	3.5	1.4	0.2
2	4.9	3.0	1.4	0.2
3	4.7	3.2	1.3	0.2
4	4.6	3.1	1.5	0.2
5	5.0	3.6	1.4	0.2
6	5.4	3.9	1.7	0.4
7	4.6	3.4	1.4	0.3
8	5.0	3.4	1.5	0.2
9	4.4	2.9	1.4	0.2
10	4.9	3.1	1.5	0.1



Supervised and Unsupervised Learning

- Supervised Learning
 - For every example in the data there is always a predefined outcome
 - Models the relations between a set of descriptive features and a target (Fits data to a function)
 - 2 groups of problems:
 - Classification
 - Regression

Supervised Learning

- **Classification**

- Predicts which class a given sample of data (sample of descriptive features) is part of (**discrete value**).

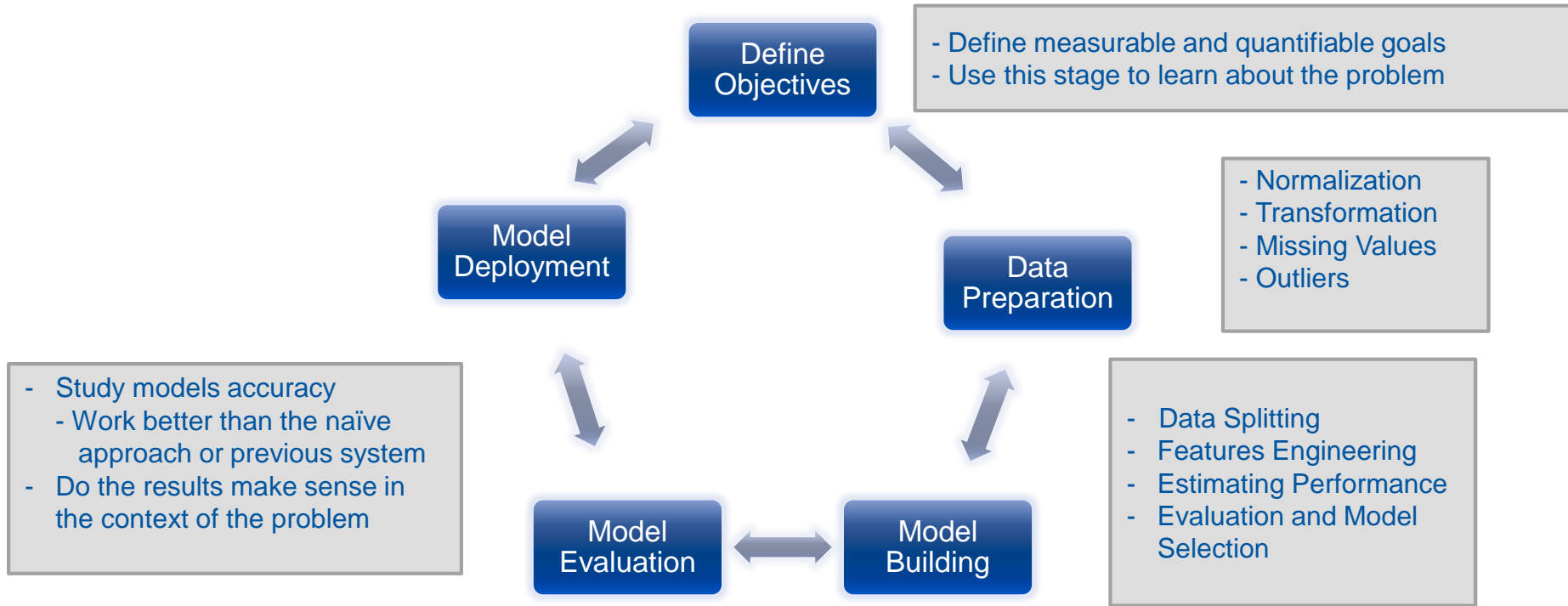
	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3.0	1.4	0.2	setosa
3	4.7	3.2	1.3	0.2	setosa
4	4.6	3.1	1.5	0.2	setosa
5	5.0	3.6	1.4	0.2	setosa
6	5.4	3.9	1.7	0.4	setosa
7	4.6	3.4	1.4	0.3	setosa
8	5.0	3.4	1.5	0.2	setosa
9	4.4	2.9	1.4	0.2	setosa
10	4.9	3.1	1.5	0.1	setosa

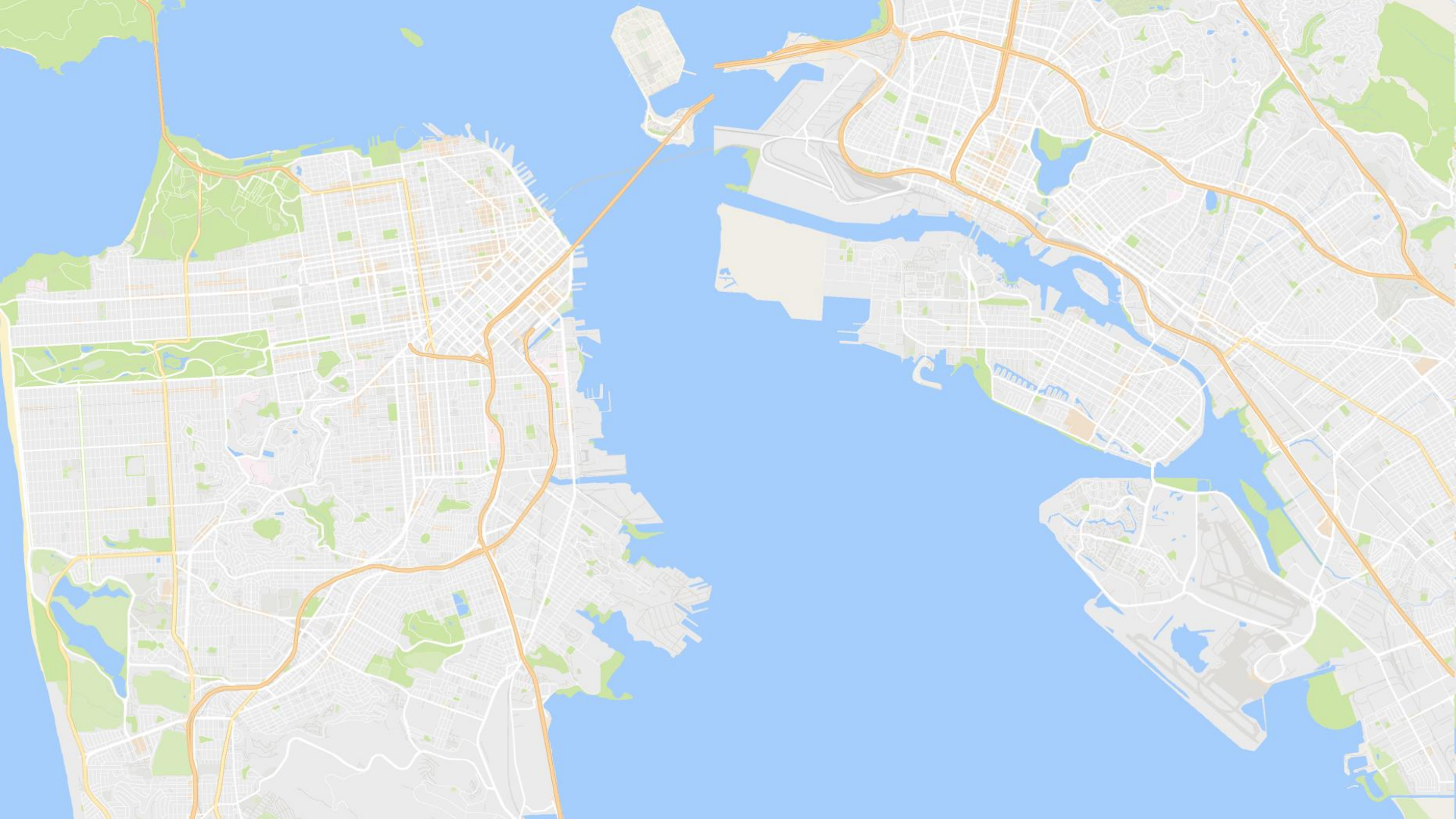
- **Regression**

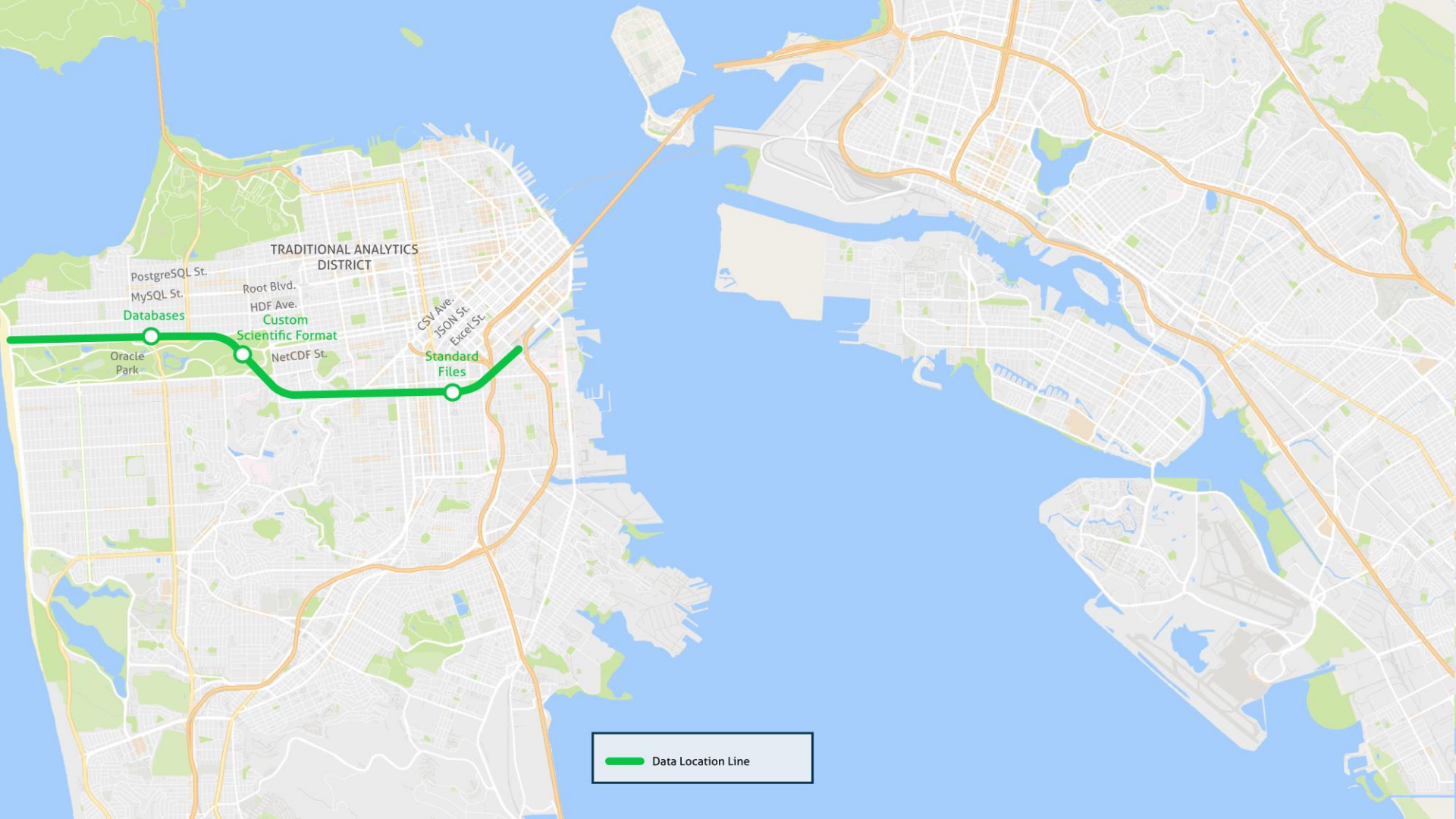
- Predicts continuous values.



Machine Learning as a Process







TRADITIONAL ANALYTICS DISTRICT

PostgreSQL St.
MySQL St.

Databases

Root Blvd.
HDF Ave.

Custom
Scientific Format

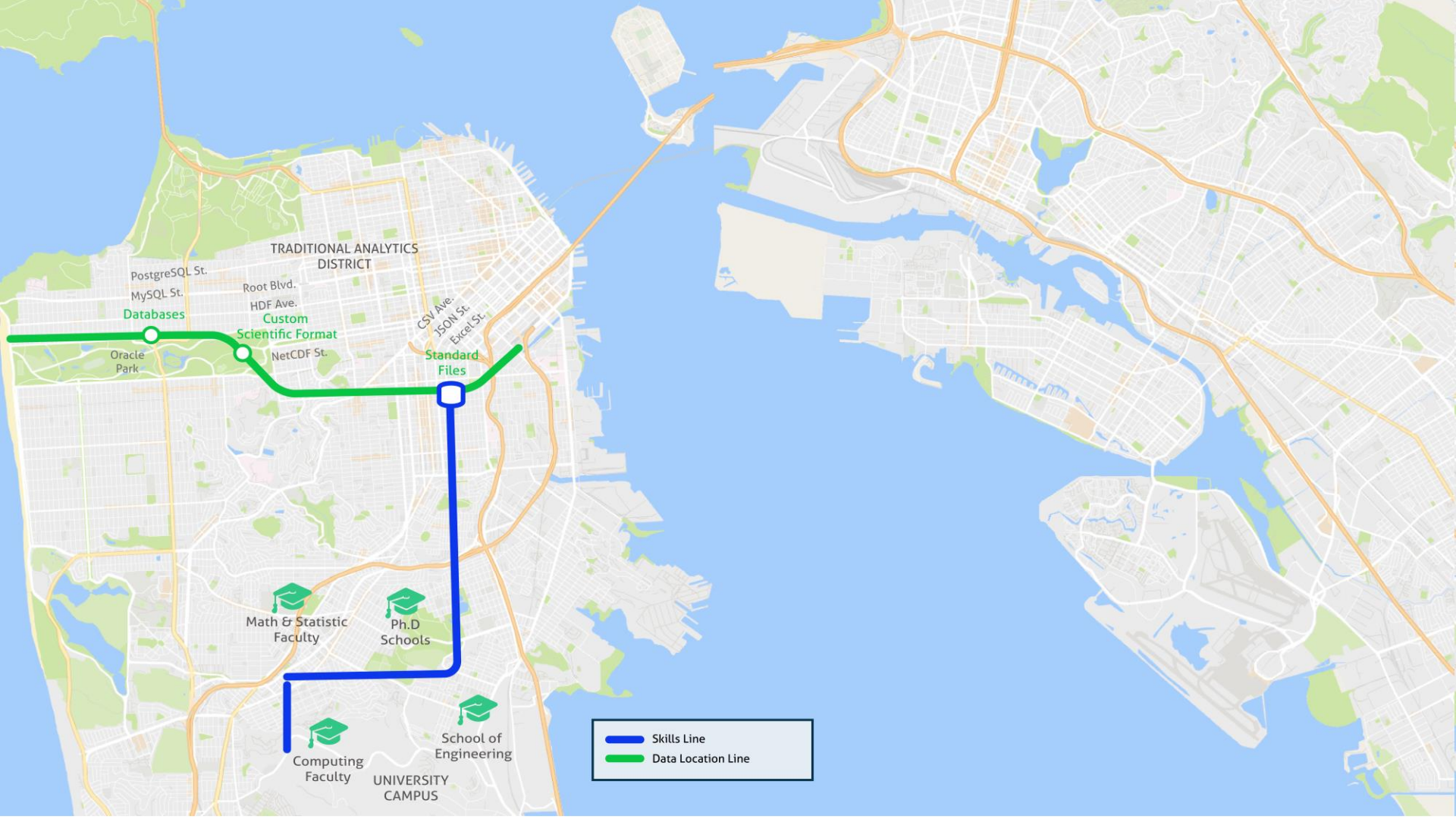
NetCDF St.

CSV Ave.
JSON St.
Excel St.

Standard
Files

Oracle
Park

 Data Location Line



TRADITIONAL ANALYTICS DISTRICT

PostgreSQL St.
MySQL St.
Databases

Root Blvd.
HDF Ave.
Custom
Scientific Format

CSV Ave.
JSON St.
Excel St.

Oracle Park

Standard Files

Math & Statistic Faculty

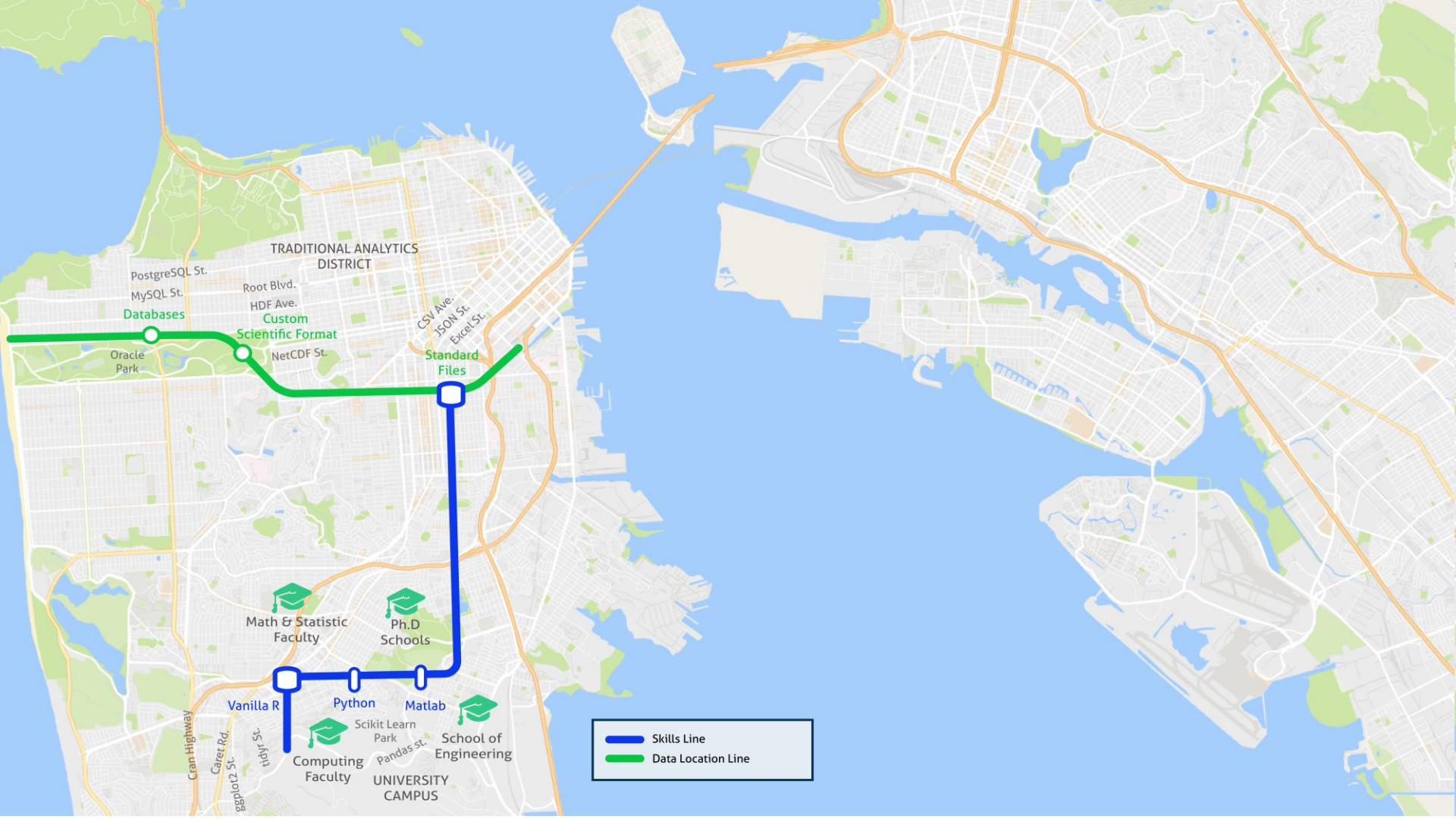
Ph.D Schools

Computing Faculty

School of Engineering

UNIVERSITY CAMPUS

Skills Line
Data Location Line



TRADITIONAL ANALYTICS DISTRICT

PostgreSQL St.
MySQL St.
Databases

Root Blvd.
HDF Ave.
Custom Scientific Format
NetCDF St.

CSV Ave.
JSON St.
Excel St.
Standard Files

Oracle Park


Math & Statistic Faculty


Ph.D Schools

Vanilla R

Python

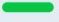
Matlab

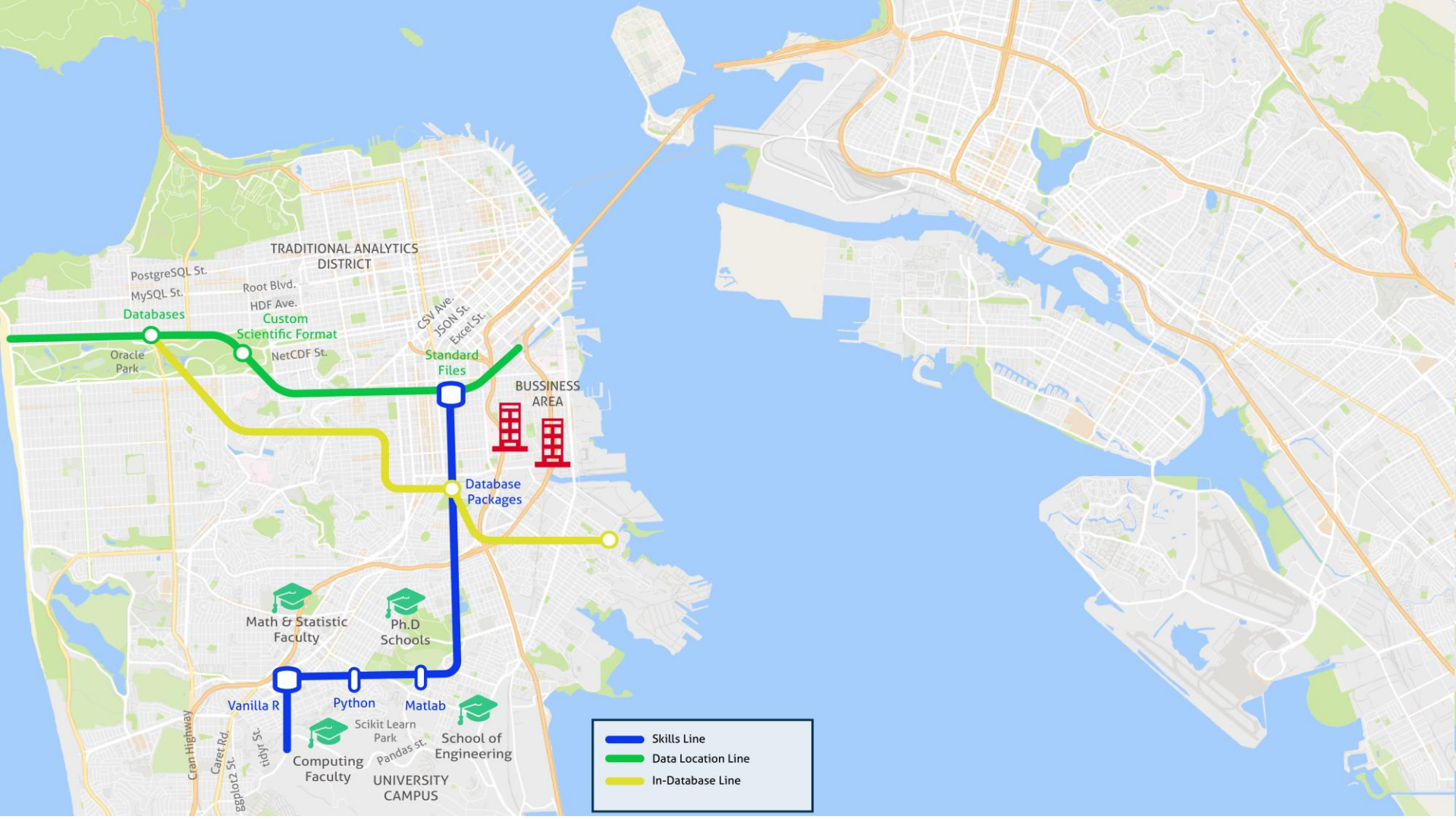

Computing Faculty

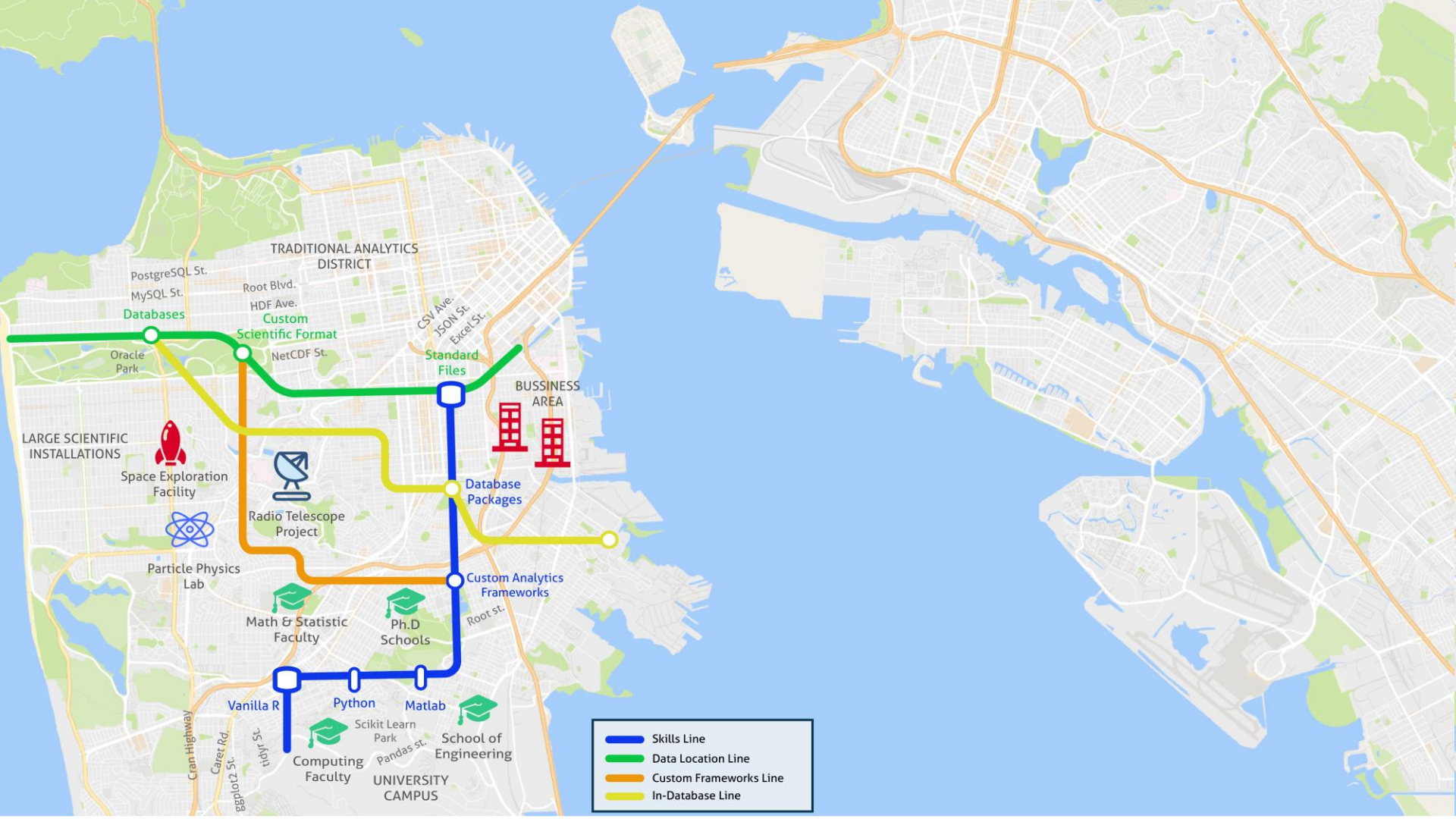
Scikit Learn Park
Pandas St.


School of Engineering

UNIVERSITY CAMPUS

	Skills Line
	Data Location Line





TRADITIONAL ANALYTICS DISTRICT

BUSINESS AREA

LARGE SCIENTIFIC INSTALLATIONS

UNIVERSITY CAMPUS

- Skills Line
- Data Location Line
- Custom Frameworks Line
- In-Database Line

PostgreSQL St.
MySQL St.
Databases

Root Blvd.
HDF Ave.
Custom Scientific Format

CSV Ave.
150th St.
Excel St.
Standard Files

Space Exploration Facility

Radio Telescope Project

Particle Physics Lab

Math & Statistic Faculty

Ph.D Schools

Computing Faculty

School of Engineering

Database Packages

Custom Analytics Frameworks

Python

Matlab

Scikit Learn Park

Pandas St.

Vanilla R

Crain Highway

Cayer Rd.

ggg102 St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

11th St.

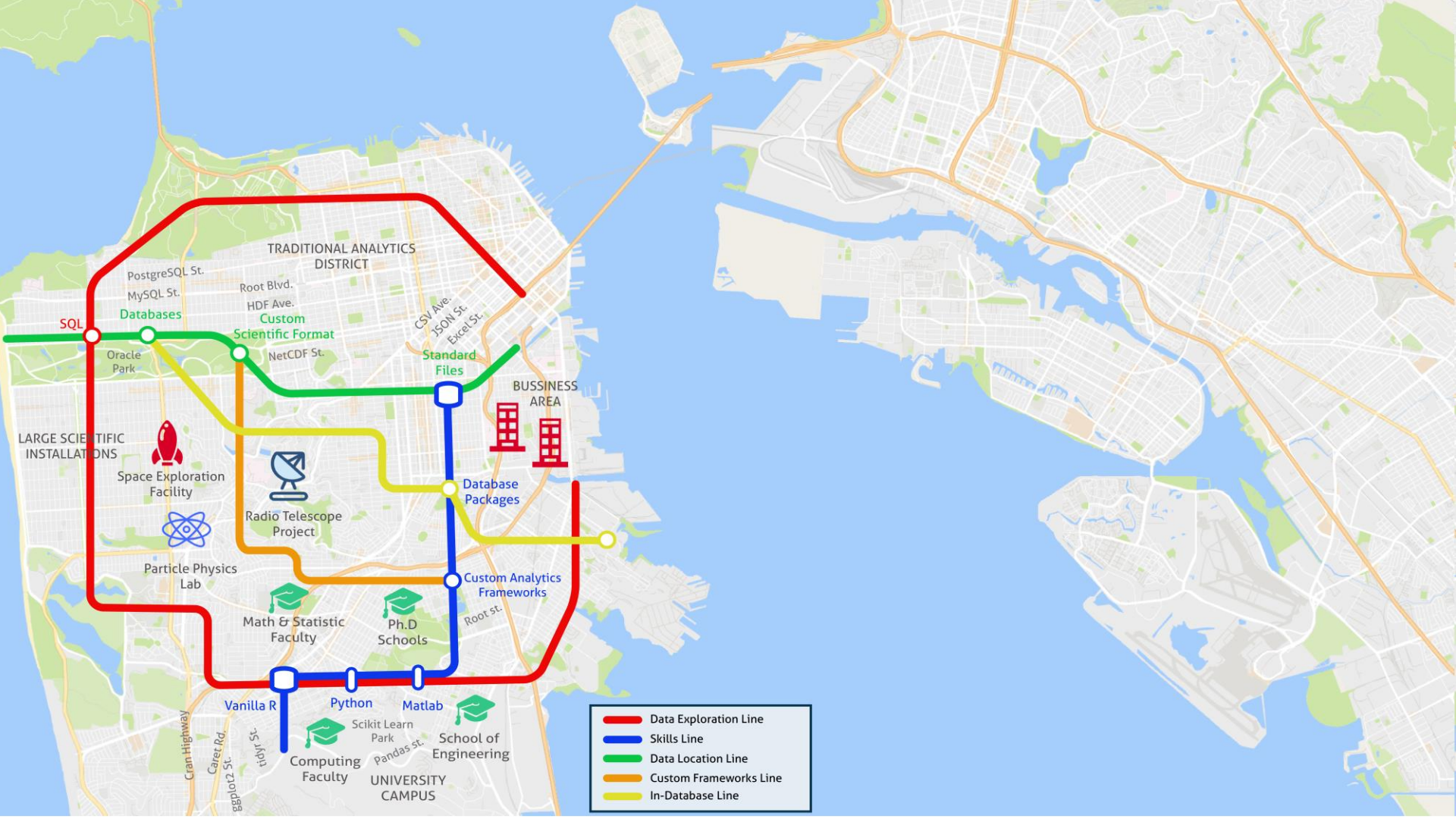
11th St.

11th St.

11th St.

11th St.

11th St.



TRADITIONAL ANALYTICS DISTRICT

BUSINESS AREA

LARGE SCIENTIFIC INSTALLATIONS

Space Exploration Facility

Particle Physics Lab

Math & Statistic Faculty

Ph.D Schools

Vanilla R
Python
Matlab
Scikit Learn
Pandas
School of Engineering
UNIVERSITY CAMPUS

Computing Faculty

Radio Telescope Project

Database Packages

Custom Analytics Frameworks

Standard Files

CSV Ave.
JSON St.
Excel St.

Root Blvd.
HDF Ave.
Custom Scientific Format

NetCDF St.

PostgreSQL St.
MySQL St.
Databases

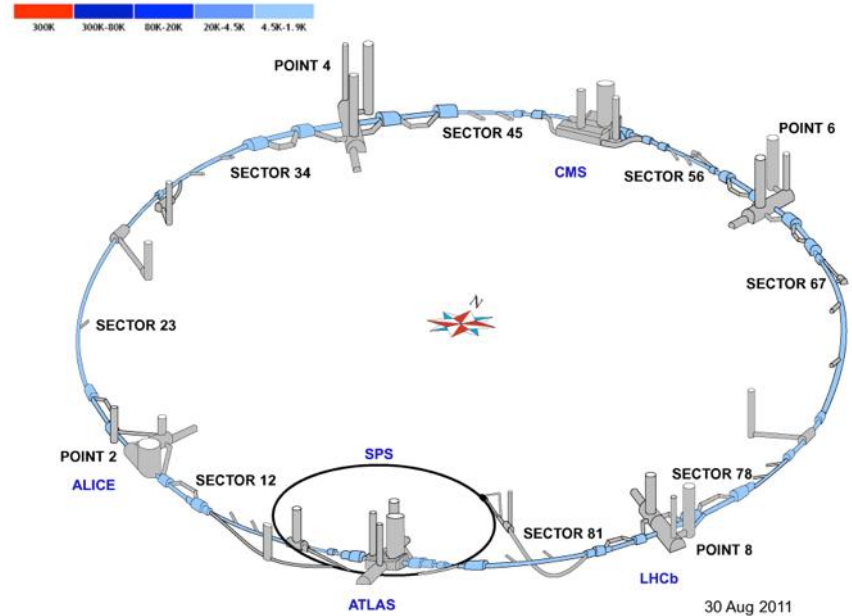
SQL

- Data Exploration Line
- Skills Line
- Data Location Line
- Custom Frameworks Line
- In-Database Line

Largest Cryogenics Installation

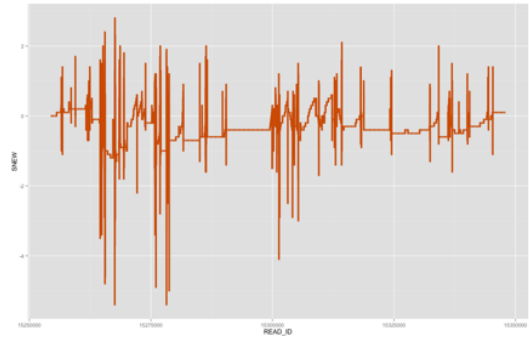
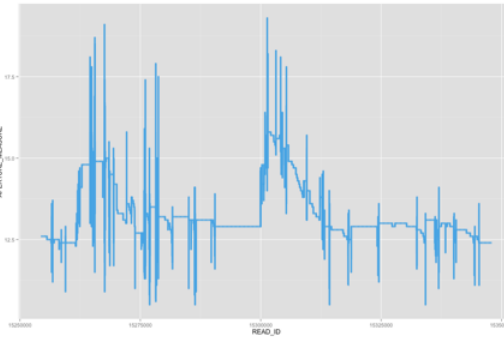
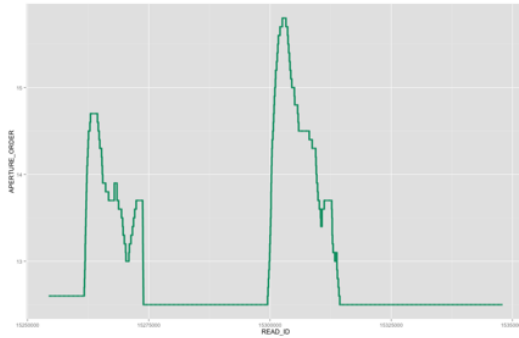
- 50k I/O, 11k actuators, ~5k control loops
- Control:
 - ~100 PLCs (Siemens, Schneider)
 - ~40 FECs (industrial PCs)
- Supervision: 26 SCADA servers

Instrument/Actuators	Total
Temperature [1.6 – 300 K]	10361
Pressure [0 – 20 bar]	2300
Level	923
Flow	2633
Control valves	3692
On/Off valves	1835
Manual valves	1916
Virtual flow meters	325
Controllers (PID)	4833



Use Case: Faulty Cryogenics Valve Detection

- What is the objective?
 - Predict faulty valves before they actually fail
- How?
 - Valve receive an aperture order value (**aperture order**)
 - Effective aperture realized by the valve (**aperture measured**)
 - Analyzing the difference between both (**S = aperture order - aperture measured**)



Faulty Cryogenics Valve Detection with R

- Signals used:
 - **S** = aperture order - aperture measured
- Features extractions based on **S**

- Variance
- Percentile 99.9
- Rope distance – R(S)
- Noise Band – B(S) (Pxx be the power spectrum of the signal S, from 0 to 0.5Hz, where S has been previously mean-centred).

$$R(S) = \frac{1}{N} \sum_{i=2}^N |S(i) - S(i-1)|$$

$$B(S) = \frac{\left| \sum_{k=1}^{N_{fft}/2} P_{xx}(k) \right|^2}{\sum_{k=1}^{N_{fft}/2} P_{xx}(k)}$$

- Automatic Faulty Valves Detection System
 - SVM - Support Vector Machine



Faulty Cryogenics Valve Detection with R

```
1 #Libraries to be used
2 library(e1071)
3 library(ggplot2)
4 library(plyr)
5 library(bspec)
6 library(randomForest)
7
8 #Reading the dataset from a CVS file (about 200Mb - reduced version)
9 CV910_dataset <- read.csv("~/Documents/IT-DB-DBB/Projects/EN Department/POST MORTEM/Data/CV910_dataset.csv")
10
11 #Calculating the diference signal
12 CV910_dataset$s<-CV910_dataset$order-CV910_dataset$measure
13
14
15 #Calculating the features (agreedated values)
16 #Sequential Calculation
17 valve.features <- ddpoly(CV910_dataset, c("cycle_id"), summarize,
18   var = var(s),
19   max = max(s),
20   min = min (s),
21   rope_dist = sum(abs(diff(s)))/length(s),
22   bs=(sum(welchPSD(s-mean(s),seglength=512,two.sided = TRUE)$power^2)/(2*512*sum(welchPSD(s-mean(s),seglength=512,two.sided = TRUE)$power^2))),
23   status = unique(status)
24 )
25
26 #Training the model - pretty fast training
27 svm_mod <- svm(status ~ var + max + min + rope_dist + bs, valve.features,na.action=na.omit)
28 rf_mod <- randomForest(status ~ var + max + min + rope_dist + bs, valve.features, na.action=na.omit, ntree=50, norm.votes=FALSE)
```

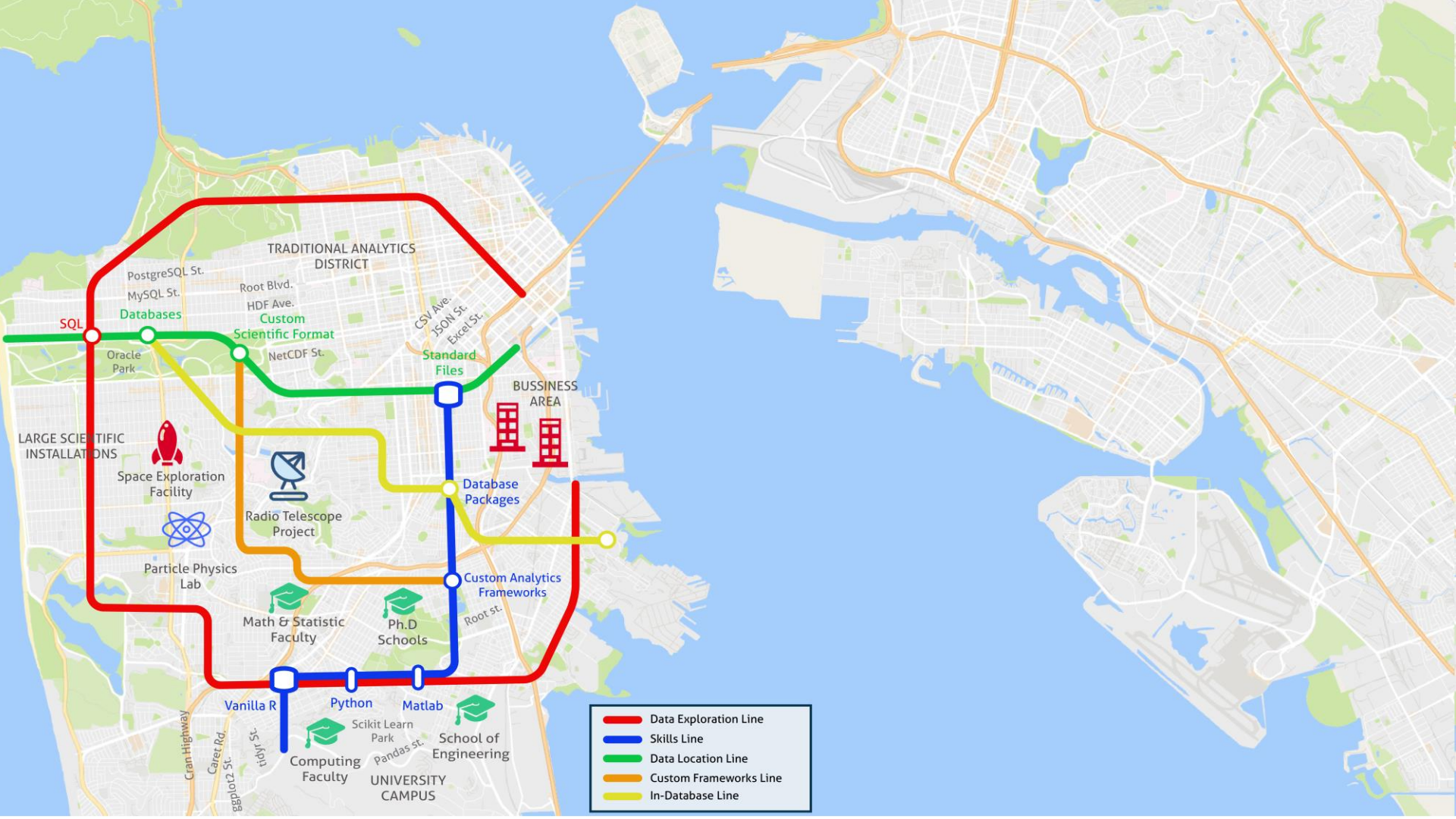
Global Environment

Data

CV910_dataset	5053287 obs. of 7 variables
valve.features	54 obs. of 7 variables

Our experience

- Excellent for prototype potential solution or validate idea
 - Fast development using standard CRAN packages such as CARET etc.
 - Large number of models and statistic functions (+7500 packages) covering a wide range of fields
 - Data Exploration
- Use the existing skills
 - R is widely use in the domain
- **Move the data is very expensive**
 - The data need to be extracted from DB and generate files CSV
 - SQL, Java API, Custom Extraction Applications (Timber)
- **Hard to deploy models in production and scale the solutions as the data grown**
 - Data limited by memory size.
 - Few packages but very limited scalability - the models themselves do not scale
 - Foreach, Snow, Rmpi, BatchExperiments package (BatchJobs)



PostgreSQL St.
MySQL St.
Oracle Park

TRADITIONAL ANALYTICS DISTRICT

Root Blvd.
HDF Ave.
Custom Scientific Format
NetCDF St.

CSV Ave.
JSON St.
Excel St.

Standard Files

BUSINESS AREA

Database Packages

Custom Analytics Frameworks

Root St.

LARGE SCIENTIFIC INSTALLATIONS

Space Exploration Facility

Radio Telescope Project

Particle Physics Lab

Math & Statistic Faculty

Ph.D Schools

UNIVERSITY CAMPUS

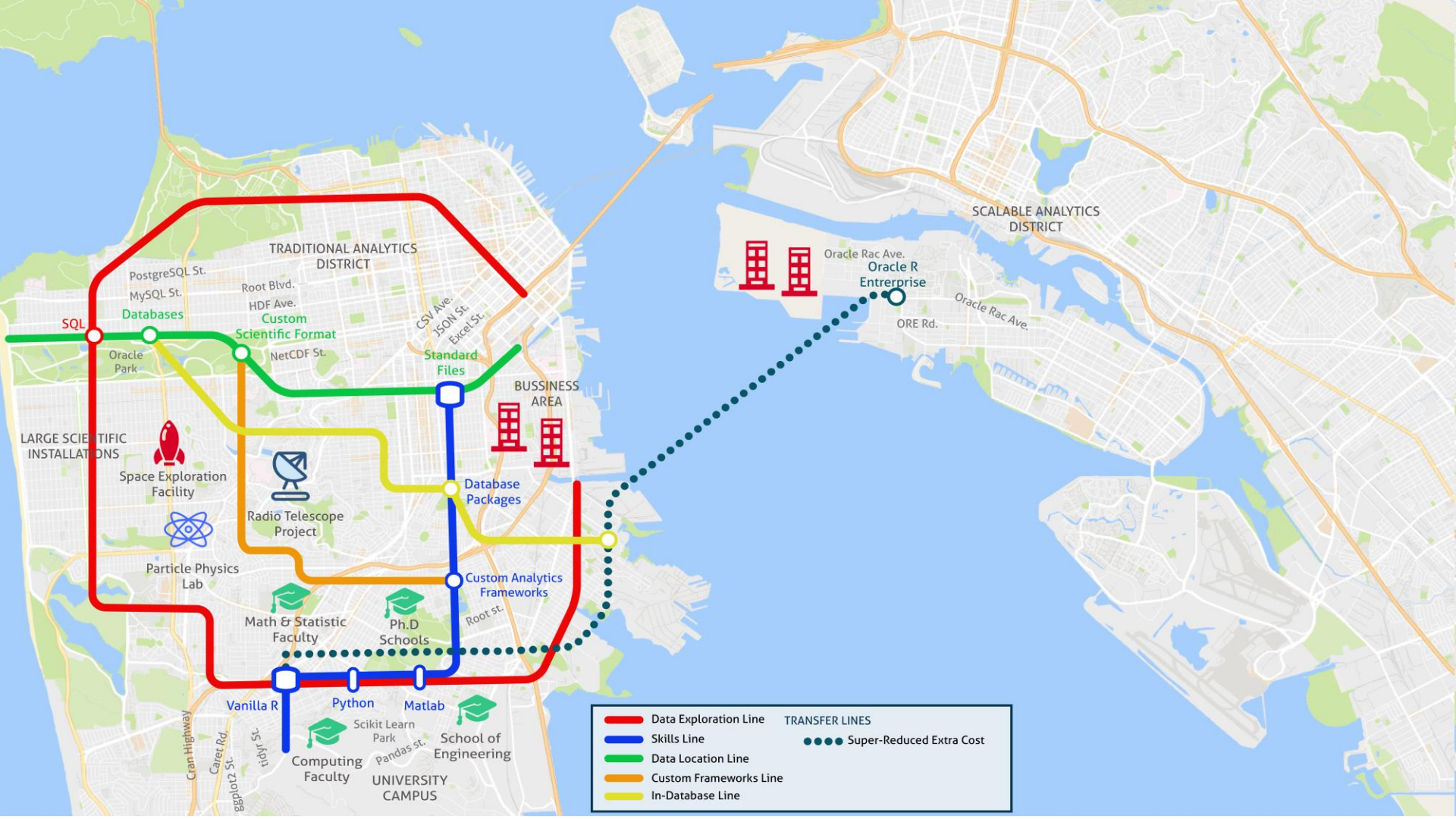
Python
Matlab
Scikit Learn Park
Pandas St.
School of Engineering

Computing Faculty

Vanilla R

Cran Highway
Carter Rd.
ggg1012 St.
tidy St.

- Data Exploration Line
- Skills Line
- Data Location Line
- Custom Frameworks Line
- In-Database Line



TRADITIONAL ANALYTICS DISTRICT

SCALABLE ANALYTICS DISTRICT

LARGE SCIENTIFIC INSTALLATIONS

BUSINESS AREA

UNIVERSITY CAMPUS

Oracle R Enterprise

—	Data Exploration Line	●●●●	TRANSFER LINES
—	Skills Line		
—	Data Location Line		
—	Custom Frameworks Line		
—	In-Database Line		

PostgreSQL St.
MySQL St.
Databases
Root Blvd.
HDF Ave.
Custom Scientific Format
NetCDF St.

CSV Ave.
JSON St.
Excel St.
Standard Files

Oracle Rac Ave.
ORE Rd.
Oracle Rac Ave.

Space Exploration Facility
Radio Telescope Project
Particle Physics Lab

Math & Statistic Faculty
Ph.D Schools
Computing Faculty

Database Packages

Custom Analytics Frameworks

Python
Matlab
School of Engineering
Scikit Learn Park
Pandas St.

Vanilla R

Crain Highway
Caret Rd.
Aggplot2 St.
tidy St.

SQL

Oracle Park

Math & Statistic Faculty

Computing Faculty

School of Engineering

Scikit Learn Park

Pandas St.

Vanilla R

Python

Matlab

School of Engineering

Scikit Learn Park

Pandas St.

Math & Statistic Faculty

Ph.D Schools

Vanilla R

Python

Matlab

School of Engineering

Scikit Learn Park

Pandas St.

Vanilla R

Python

Matlab

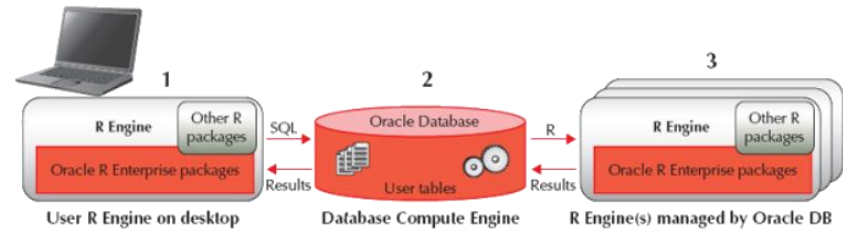
School of Engineering

Scikit Learn Park

Pandas St.

Why ORE? - ORE benefits

- A database-centric environment for analytical processes in R
 - Allows to use the database server to run R scripts (scalability & performance)
 - Eliminate memory constraint of client R engine



- **Transparency Layer**
 - Transparently analyze and use data in Oracle Database through R
 - Tables as R native data frames
- Enables users to take advantage of **data-parallel** and **task-parallel** execution through Oracle Database

Cryo Valves — Parallel Features Extraction in ORE

```
4 #Function to calculate and extract the features
5 features <- function(dat) {
6   #Load R Libraries
7   library(bspec)
8
9   #Calculate the signal to work with
10  s<-dat$APERTURE_ORDER-dat$APERTURE_MEASURE
11
12  #Features Calculations
13  valve = unique(dat$VALVE)
14  cycle = unique(dat$CYCLE_NUMBER)
15  status = unique(dat$STATUS)
16
17  var=var(s)
18  max=max(s)
19  min=min(s)
20
21  rope_dist=sum(abs(diff(s)))/length(s)
22  pxx<-welchPSD(s-mean(s),seglength=512,two.sided = TRUE)$power
23  bs=(sum(pxx)^2)/(2*512*sum(pxx^2))
24
25  #Return the features
26  data.frame(Cycle = cycle,
27            Valve = valve,
28            Status = status,
29            Var=var,
30            Max=max,
31            Min=min,
32            Rope_dist = rope_dist,
33            Bs = bs)
34
35 }
```

Instrument/Actuators	Total
Temperature [1.6 – 300 K]	10361
Pressure [0 – 20 bar]	2300
Level	923
Flow	2633
Control valves	3692
On/Off valves	1835
Manual valves	1916
Virtual flow meters	325
Controllers (PID)	4833

93600 points per cycle (about 24 hours)



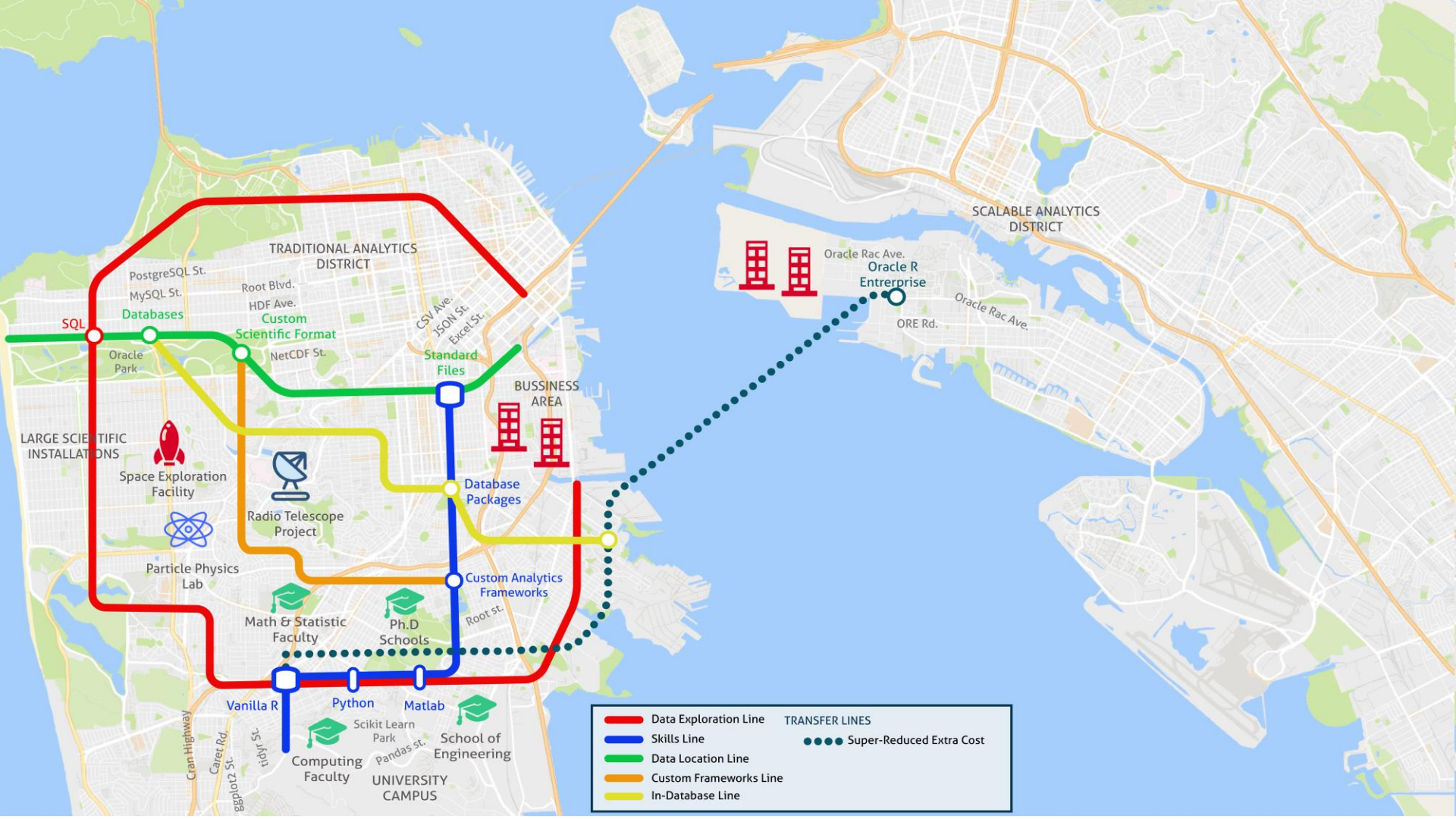
Cryo Valves — Parallel Features Extraction in ORE

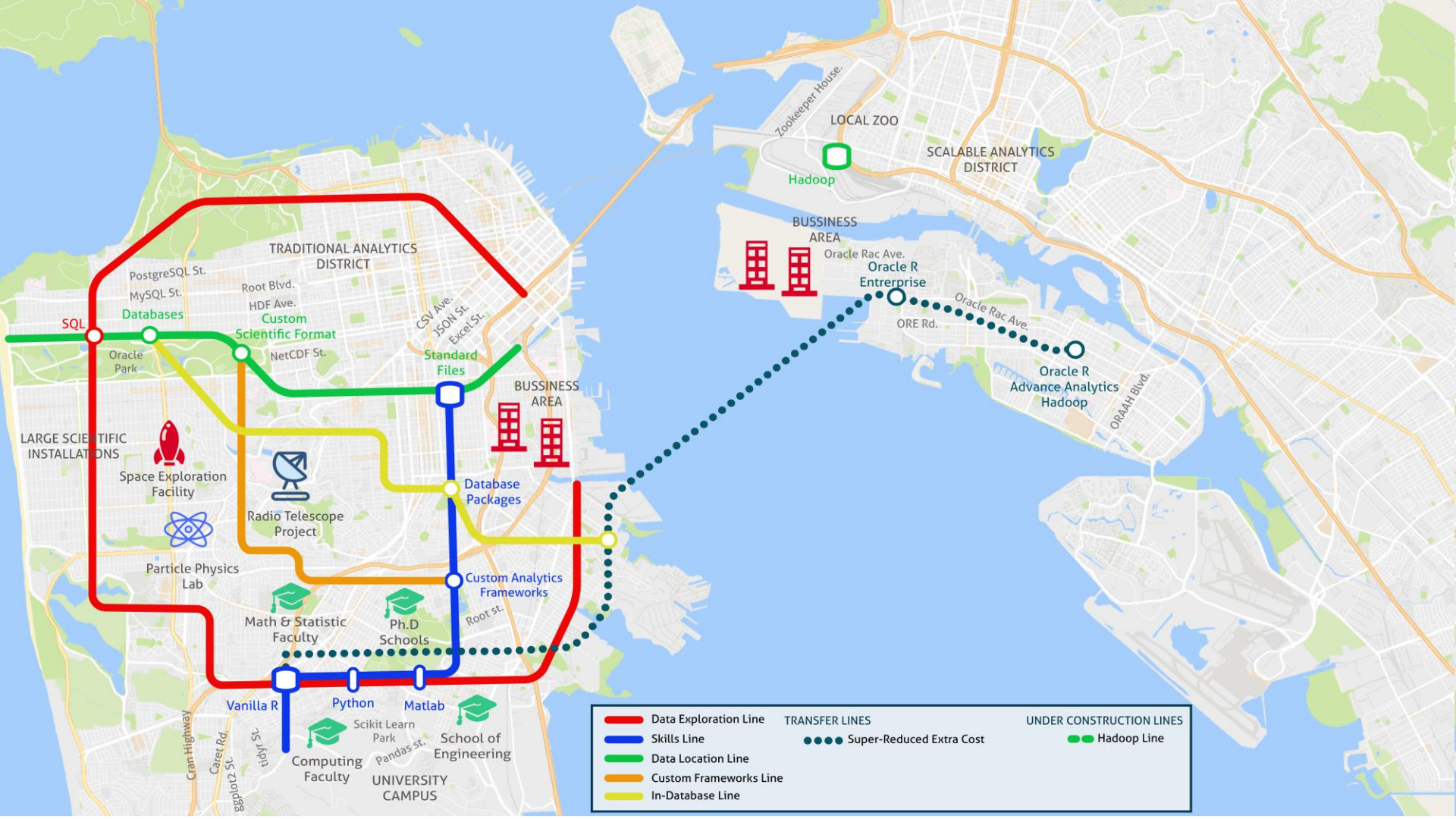
```
37 #Create in-db partitions and process them in parallel (VALVE_READS_CL_TRA is the name of the table in our db)
38 #Lazy evaluation!!!
39 ore.valve.features <- ore.groupApply(VALVE_READS_CL_TRA [,c('READ_ID', 'CYCLE_NUMBER', 'VALVE', 'APERTURE_ORDER', 'APERTURE_MEASURE', 'STATUS')],
40                                     INDEX=VALVE_READS_CL_TRA$CYCLE_NUMBER,
41                                     features,
42                                     #Return value signature
43                                     FUN.VALUE =
44                                     data.frame(Cycle = numeric(),
45                                               Valve = character(),
46                                               Status = character(),
47                                               Var = numeric(),
48                                               Max = numeric(),
49                                               Min = numeric(),
50                                               Rope_dist = numeric(),
51                                               Bs = numeric())
52                                     ,parallel=TRUE)
53
54 #Order the output
55 row.names(ore.valve.features)=ore.valve.features$Cycle
56
57 #Pulling the fatures from DB for local processing
58 valve.features<-ore.pull(ore.valve.features)
59
60 #Training the models locally
61 rf_mod <- svm(Status ~ Var + Max + Min + Rope_dist + Bs, valve.features,na.action=na.omit)
```

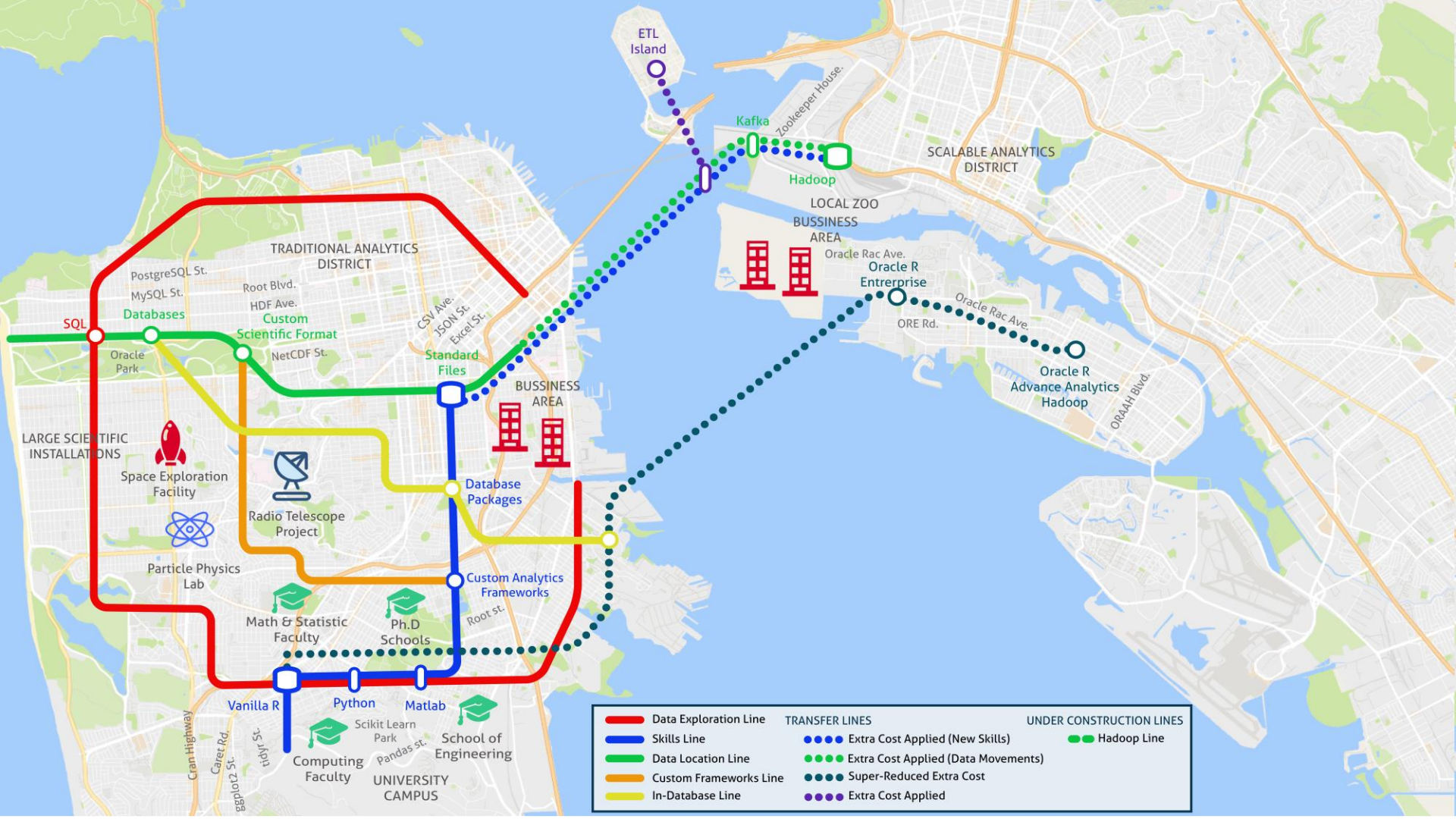


Our experience

- No need to move data
 - It is faster to process it in-DB with ORE using the appropriate degree of parallelism
- DB nodes already prepared for the workload
 - Simplifies the infrastructure
- Write/adapt R code is straight forward
 - Thanks to transparency layer and embedded R execution
 - Tables and Views as R dataframes
- **Still problems scaling**
 - Scalability determined by RAC installations
- **Need to differentiate between production and analytics environments**
 - Risk on affecting the production environment performance by running in-database analytics
- **Analytics developments in-database**
 - Risk on data security and resources competition





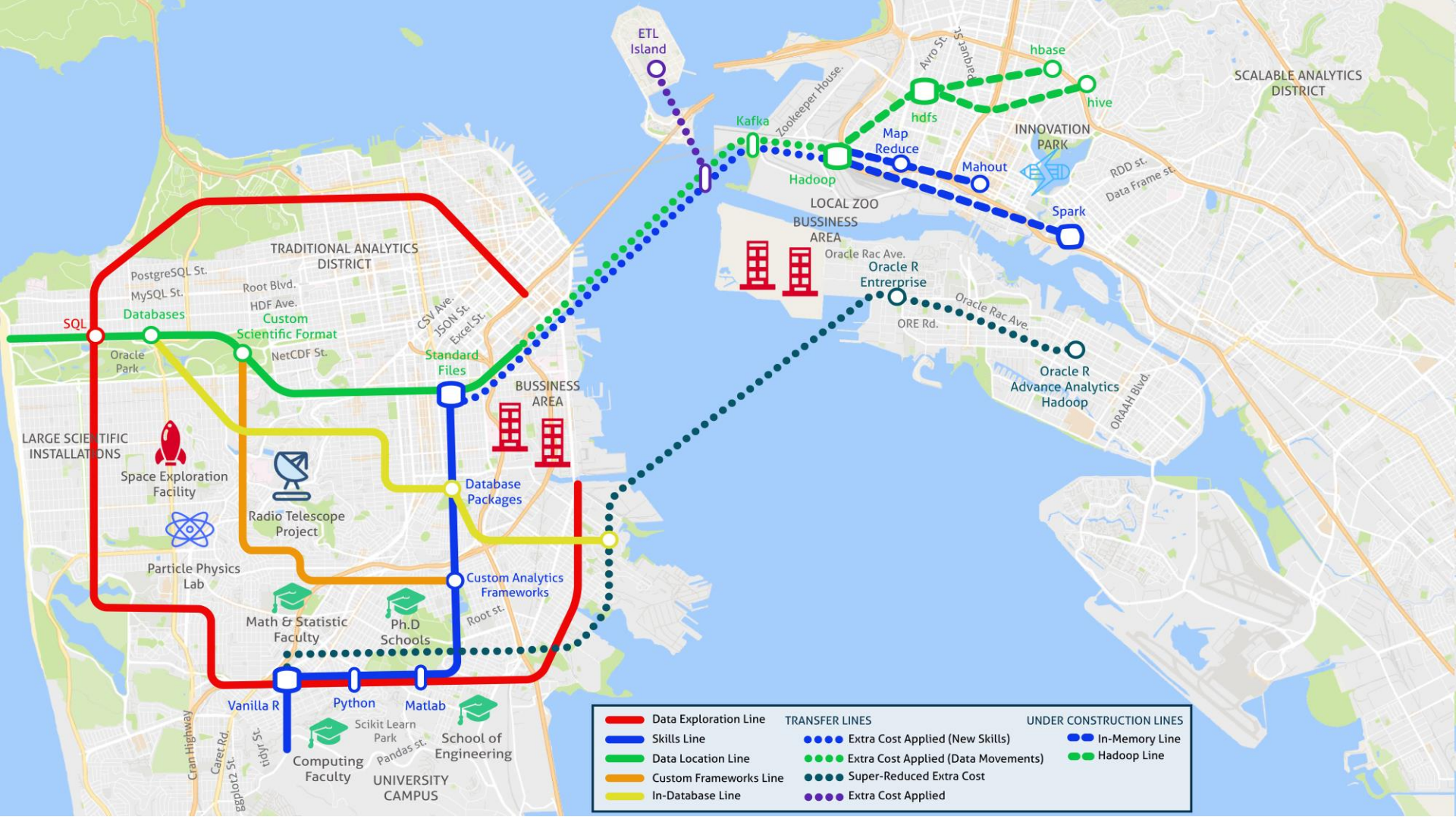


	Data Exploration Line		Extra Cost Applied (New Skills)		Under Construction Lines
	Skills Line		Extra Cost Applied (Data Movements)		Extra Cost Applied
	Data Location Line		Super-Reduced Extra Cost		
	Custom Frameworks Line				
	In-Database Line				

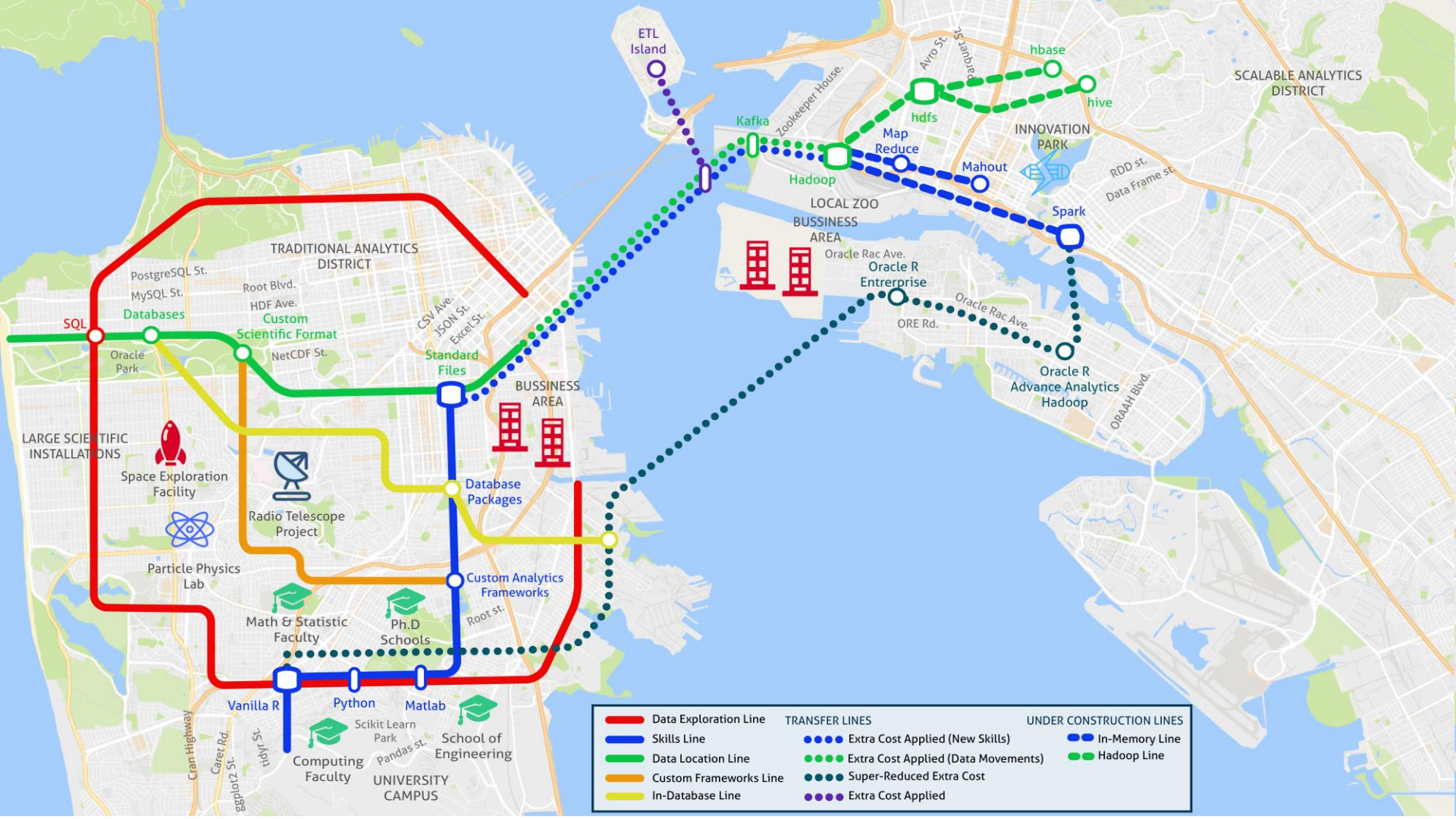


	Data Exploration Line		TRANSFER LINES Extra Cost Applied (New Skills)		UNDER CONSTRUCTION LINES In-Memory Line
	Skills Line		Extra Cost Applied (Data Movements)		Hadoop Line
	Data Location Line		Extra Cost Applied		
	Custom Frameworks Line				
	In-Database Line				

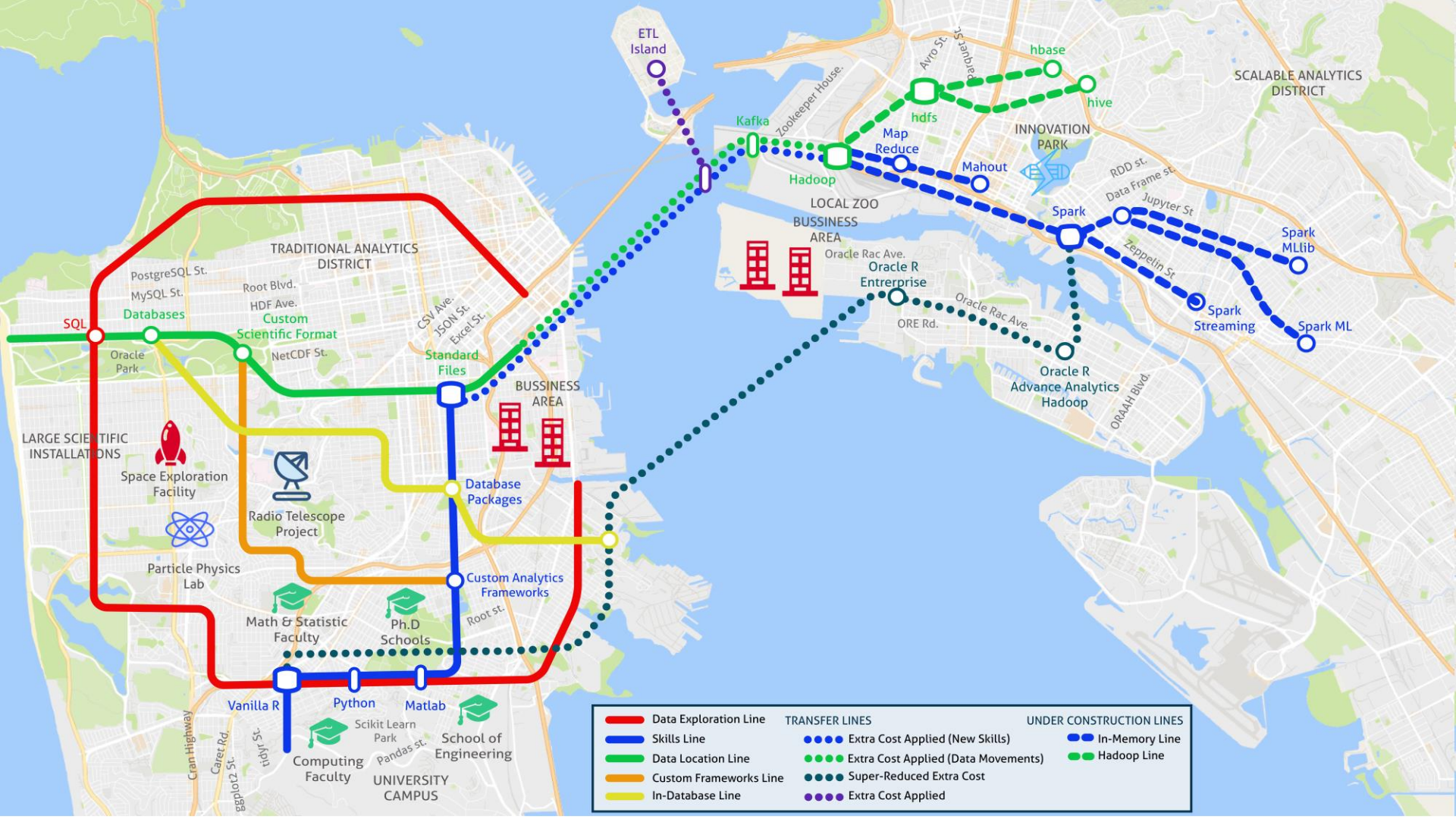




TRANSFER LINES		
— Data Exploration Line	●●●● Extra Cost Applied (New Skills)	— In-Memory Line
— Data Location Line	●●●● Extra Cost Applied (Data Movements)	— Hadoop Line
— Custom Frameworks Line	●●●● Super-Reduced Extra Cost	
— In-Database Line	●●●● Extra Cost Applied	



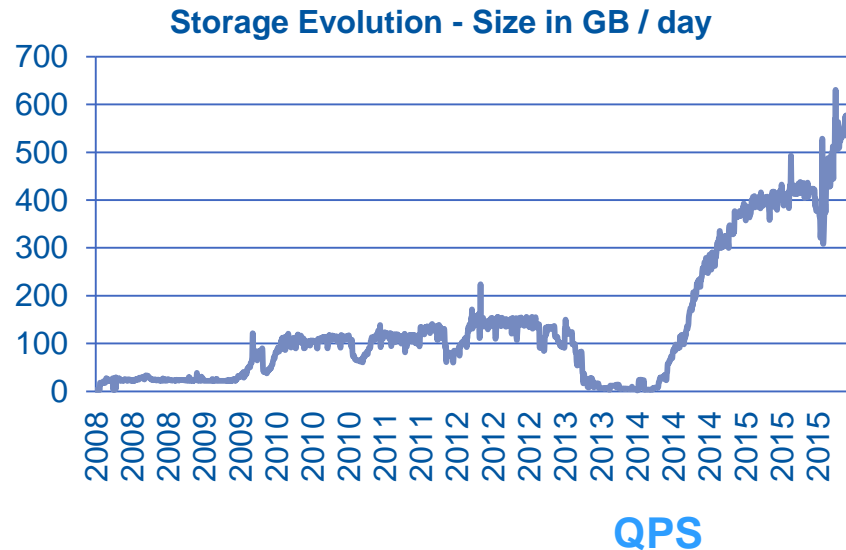
TRANSFER LINES		
— Data Exploration Line	●●●● Extra Cost Applied (New Skills)	- - - In-Memory Line
— Skills Line	●●●● Extra Cost Applied (Data Movements)	- - - Hadoop Line
— Data Location Line	●●●● Super-Reduced Extra Cost	
— Custom Frameworks Line	●●●● Extra Cost Applied	
— In-Database Line		



TRANSFER LINES					
	Data Exploration Line		Extra Cost Applied (New Skills)		In-Memory Line
	Skills Line		Extra Cost Applied (Data Movements)		Hadoop Line
	Data Location Line		Super-Reduced Extra Cost		
	Custom Frameworks Line		Extra Cost Applied		
	In-Database Line				

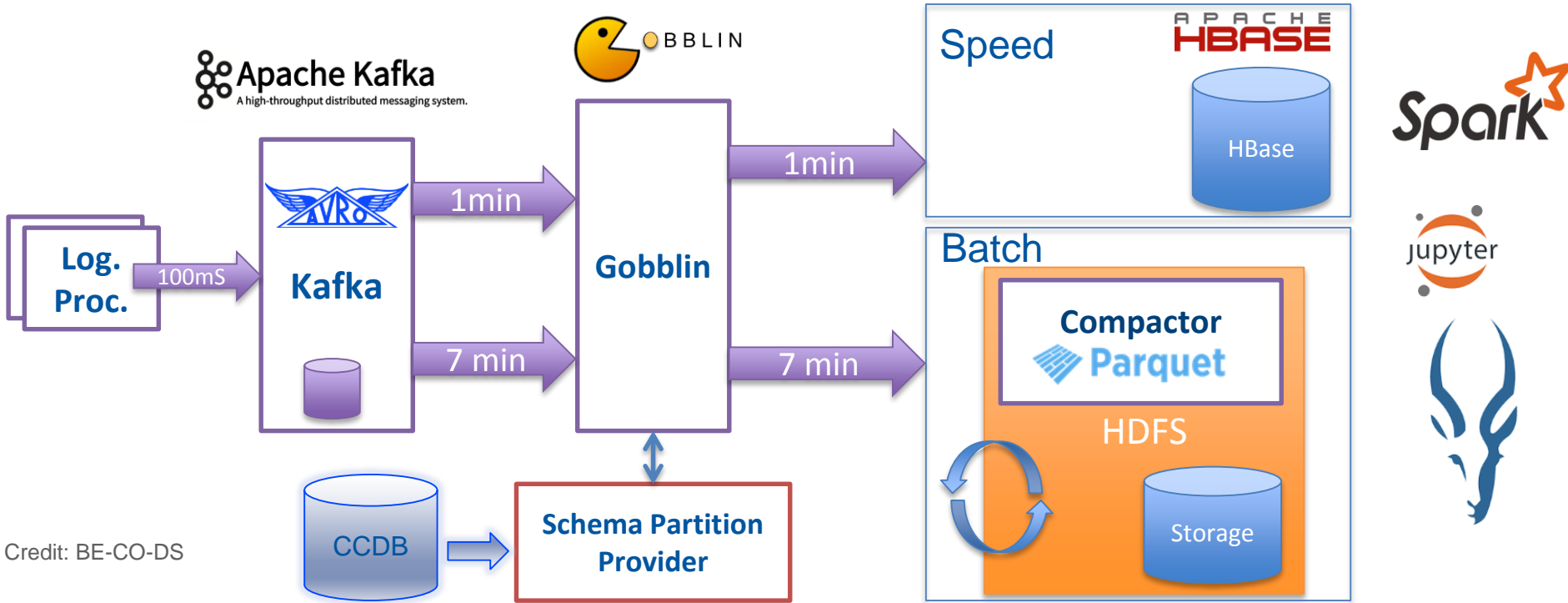
CERN Accelerator Logging Service 2.0

- New Landscape bring new challenges
 - Better Performance on bigger datasets
 - Big Data queries: Impala, Spark SQL
 - Leverage analytics capabilities
 - Spark Analytics: Python, ML, R
 - More heterogeneous data access models



Credit: BE-CO-DS

CERN Accelerator Logging Service



Credit: BE-CO-DS

Cryo Valves – ORAAH

Oracle R Advanced Analytics
for Hadoop

```
1 library(rpart)
2 library(ORCH)
3
4 #Setting up the connection to Spark installation
5 spark.connect(master="yarn-client",memory="8G",dfs.namenode="-----")
6
7 #Preparing the data on HDFS
8 train <- hdfs.put(training)
9 test <- hdfs.put(testing)
10
11 #Let focus now in model generation and training
12 #Spark MLlib algorithms available from R in ORAAH (source data can be CSVs in HDFS or HIVE tables)
13 model <- orch.ml.svm(formula = status ~ rope_dist + bs + mean + var + max ,data = train)
14
15 #Predincting on test and writing back the results to HDFS
16 pred <- predict(model, newdata = test)
17 hdfs.write(pred, outPath = "Prediction")
```


Our experience

- No need to move data – the analysis is done where the data is
 - Access to database or Hive tables transparently
- Memory is not a problem anymore
 - The analysis is not anymore limited by the dataset size
 - Simplifies the infrastructure – transparent use of hadoop technologies
- Write/adapt R code is straight forward
 - Same principals than ORE – no need to acquire a new set of skills
 - Background use of Spark machine learning capabilities
- **Limited functionality**
 - When using ORAAH for Machine learning in scalable way the functionality is limited to Spark Machine Learning libraries
 - No as fast pace as Spark itself – Why do we need to wait
- **Commercial VS open source**

Machine Learning with Spark

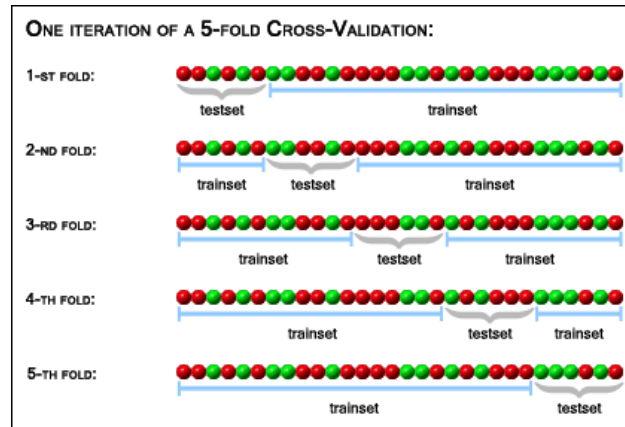


- Why Apache Spark for Machine Learning
 - No memory limitations
 - Compatibility – Scala, R, Python
 - General purpose
 - Not only machine learning also **advanced data preparation, feature engineering, parameter tuning and model selection** etc.
- New Skills required
 - RDD-based API (spark.mllib)
 - DataFrame-based API (spark.ml)
 - Pipelines Concept (CARET package does in R)
 - Cross Validation
 - Parameter tuning using parameter grid
- Fast Pace Evolution

Machine Learning with Spark

- **Cross Validation**

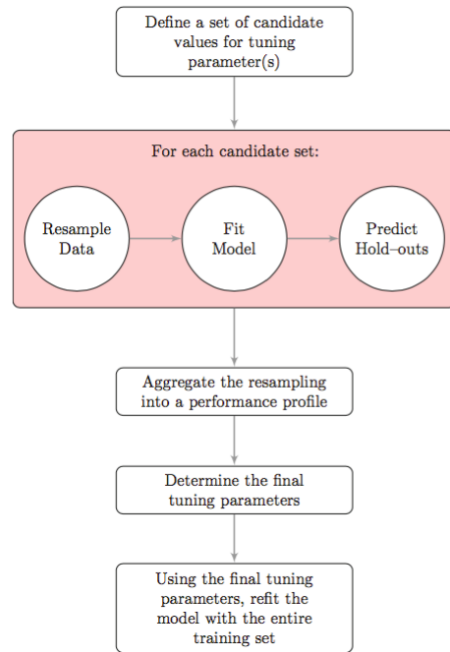
- Repeat the construction of the model on different subsets of the available training data and then evaluate the model only on data not seen during training



Machine Learning with Spark

- ## Model Tuning

- ML models have several parameter
 - there is **no** analytics formula to calculate appropriate values
- These parameters control the complexity of the model
 - bad performance
 - over-fitting
 - etc.



Cryo Valves – Spark

```
In [ ]: from pyspark.ml.classification import RandomForestClassifier

#select the features to used
features = ['var', 'max', 'min', 'rope_dist', 'bs']

#We have to shape our dataset before we can use some Spark ML algorithms with it.
#First, we need to create a new column containing a Vector with all the features.
#That vector will be the input features used by the ML for the model training.
assembler = VectorAssembler(inputCols=features, outputCol="features")

rfClassifierCV = RandomForestClassifier(labelCol="status", featuresCol="features")

#Setting the pipeline
rfPipelineCV = Pipeline(stages=[assembler, rfClassifierCV])

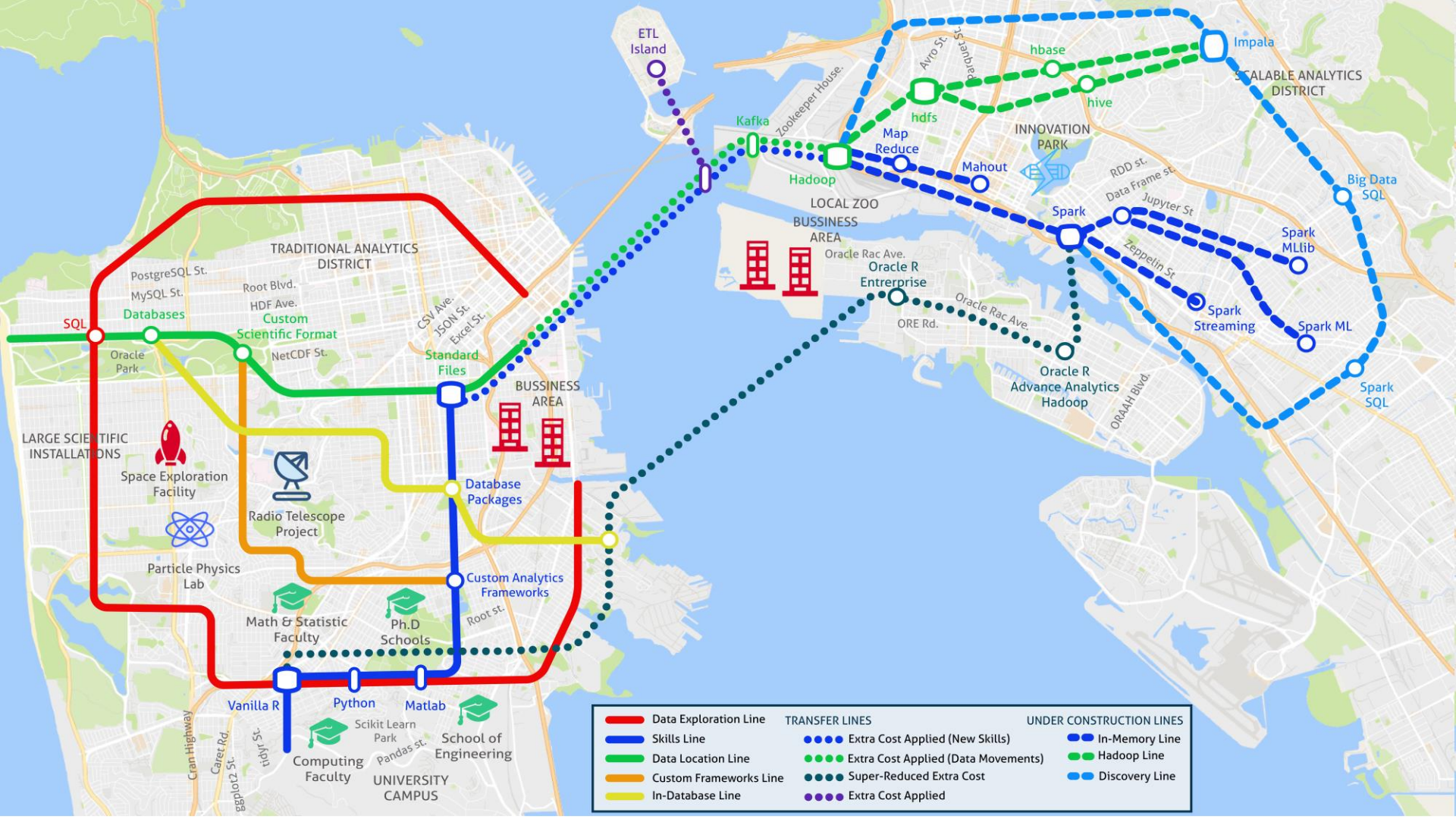
#Prepare de evaluator, CV configuration and parameter tuning
evaluatorCV = BinaryClassificationEvaluator(labelCol="indexedLabel")

#3 different models will be evaluated 5,10 and 15 number of trees
paramGrid = ParamGridBuilder().addGrid(rfClassifierCV.numTrees, [5,10,15]).build()

#Definition of a the CV pofile
crossValidator = CrossValidator(estimator=rfPipelineCV, estimatorParamMaps=paramGrid, evaluator=evaluatorCV, numFolds=3)

#Now we perform the computation all the prevoius code is lazy evaluated
rfClassifierCV = crossValidator.fit(train)
rfPredictionCV = cvModel.transform(test)
```





TRANSFER LINES		
— Data Exploration Line	●●●● Extra Cost Applied (New Skills)	- - - - In-Memory Line
— Skills Line	●●●● Extra Cost Applied (Data Movements)	- - - - Hadoop Line
— Data Location Line	●●●● Super-Reduced Extra Cost	- - - - Discovery Line
— Custom Frameworks Line	●●●● Extra Cost Applied	
— In-Database Line		

Machine Learning with TensorFlow

- Why TensorFlow for Machine Learning
 - Spark machine learning capabilities are really limited
 - Number of models
 - Customization capabilities
 - Overcome in term of performance any of the previous technologies
 - Spark is slow on training models
 - State-of-the-art algorithms available
 - Deep-learning
 - New skill need to be understood
 - Tensor concept
 - Model freedom comes with a price
 - Coding

Machine Learning with TensorFlow+Spark

- No memory limitations
 - Bigger than memory datasets treated transparently
- Parallelization
 - Tensorflow profit from Spark partitioning concepts to improve the user control over parallelization

