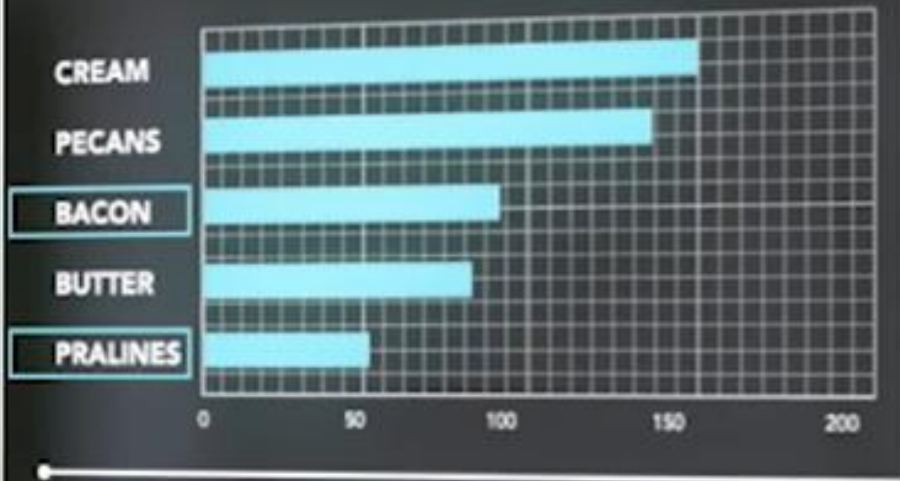


**BEST SELLER:  
PECANS & CREAM**

□ SOCIAL AFFINITY SEARCH



**SENTIMENT ANALYSIS: BACON + PRALINES**



# Microsoft R server Stefan Cronjaeger

Technical Solution Specialist Advanced Analytics  
Global Blackbelt – Germany

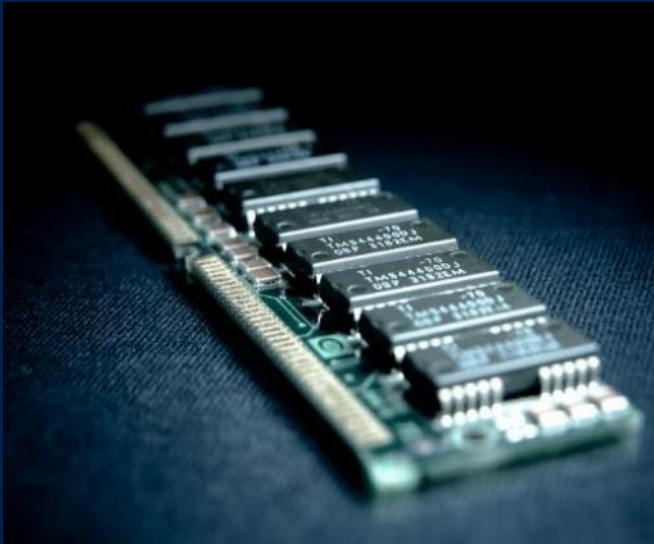
[scronj@microsoft.com](mailto:scronj@microsoft.com)

+49 151 4406 3425

# Topics

- Microsoft R Server (and demo)
- R on Hadoop
- R on SQL server
- Operationalizing R

# Enterprise use of open source R



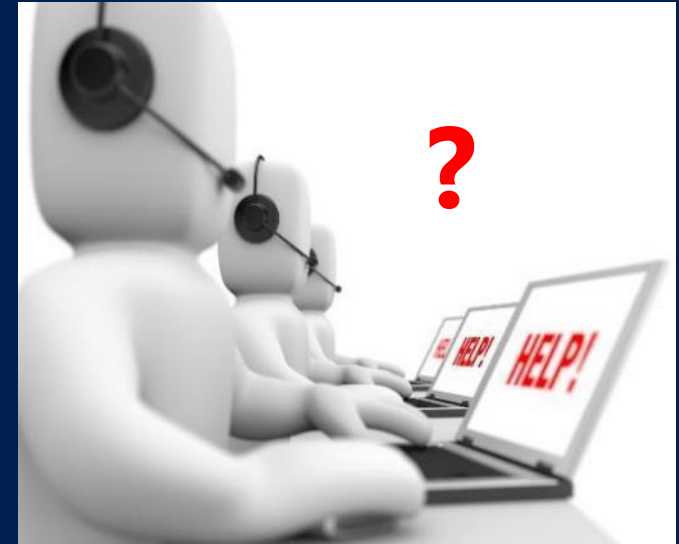
R needs data in memory  
to start a computation\*



R is mostly single  
threaded



R requires skilled  
resource to scale out  
computations across a  
cluster and needs re-  
coding for R map-  
reduce in Hadoop



Open source R is  
supported by the  
community

# CRAN, MRO, MRS Comparison



**Microsoft  
R Open**

**Microsoft  
R Server**

<b>Dataseize</b>	In-memory	In-memory	<b>In-Memory or Disk Based</b>
<b>Speed of Analysis</b>	Single threaded	Multi-threaded	<b>Multi-threaded, parallel processing 1:N servers</b>
<b>Support</b>	Community	Community	<b>Community + Commercial</b>
<b>Analytic Breadth &amp; Depth</b>	7500+ innovative analytic packages	7500+ innovative analytic packages	<b>7500+ innovative packages + commercial parallel high-speed functions</b>
<b>Licence</b>	Open Source	Open Source	<b>Commercial license. Supported release with indemnity</b>

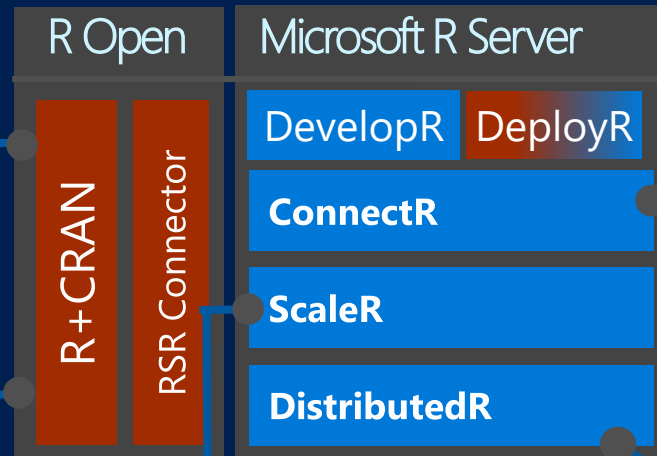
# The Microsoft R Server Platform

## R+CRAN

- Open source R interpreter
  - R 3.1.2
- Freely-available huge range of R algorithms
- 100% Compatible with existing R scripts, functions and packages

## RevoR

- Performance enhanced R interpreter
- Based on open source R
- Adds high-performance math library to speed up linear algebra functions



## ScaleR

- Ready-to-Use high-performance big data big analytics
- Fully-parallelized analytics

## ConnectR

- High-speed & direct connectors

### Available for:

- High-performance XDF
- SAS, SPSS, delimited & fixed format text data files
- Hadoop HDFS (text & XDF)
- Teradata Database & Aster
- ODBC

## DistributedR

- Distributed computing framework
- Delivers cross-platform portability

# Scale R – Parallelized Algorithms & Functions

## Data Preparation

- Data import – Delimited, Fixed, SAS, SPSS, ODBC
- Variable creation & transformation
- Recode variables
- Factor variables
- Missing value handling
- Sort, Merge, Split
- Aggregate by category (means, sums)

## Descriptive Statistics

- Min / Max, Mean, Median (approx.)
- Quantiles (approx.)
- Standard Deviation
- Variance
- Correlation
- Covariance
- Sum of Squares (cross product matrix for set variables)
- Pairwise Cross tabs
- Risk Ratio & Odds Ratio
- Cross-Tabulation of Data (standard tables & long form)
- Marginal Summaries of Cross Tabulations

## Statistical Tests

- Chi Square Test
- Kendall Rank Correlation
- Fisher's Exact Test
- Student's t-Test

## Sampling

- Subsample (observations & variables)
- Random Sampling

## Predictive Models

- Sum of Squares (cross product matrix for set variables)
- Multiple Linear Regression
- Generalized Linear Models (GLM) exponential family distributions: binomial, Gaussian, inverse Gaussian, Poisson, Tweedie. Standard link functions: cauchit, identity, log, logit, probit. User defined distributions & link functions.
- Covariance & Correlation Matrices
- Logistic Regression
- Classification & Regression Trees
- Predictions/scoring for models
- Residuals for all models

## Variable Selection

- Stepwise Regression

## Simulation

- Simulation (e.g. Monte Carlo)
- Parallel Random Number Generation

## Cluster Analysis

- K-Means

## Classification

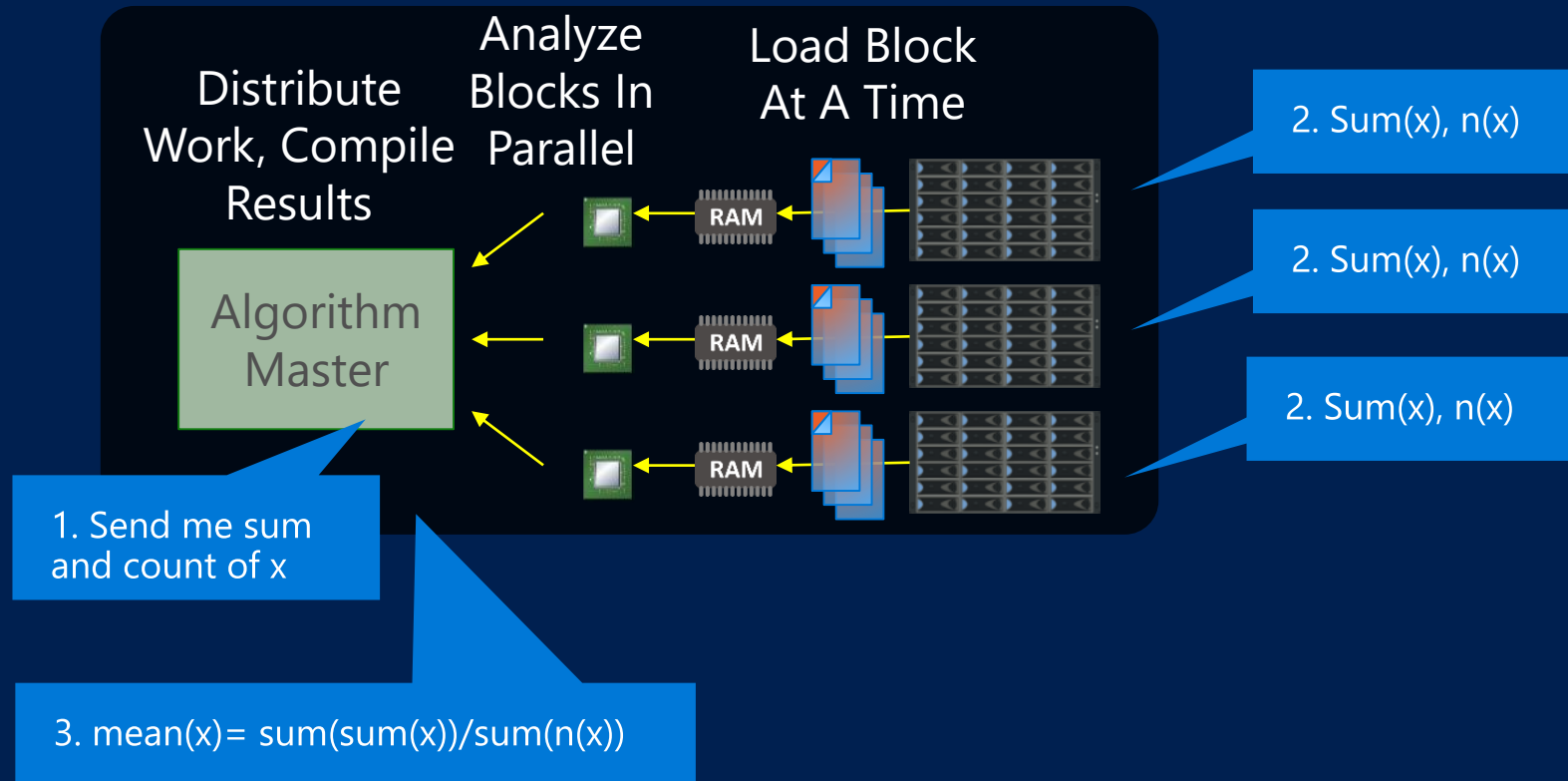
- Decision Trees
- Decision Forest
- Gradient Boosted Decision Trees
- Naïve Bayes



## Combination

- rxDataStep
- rxExec
- PEMA API

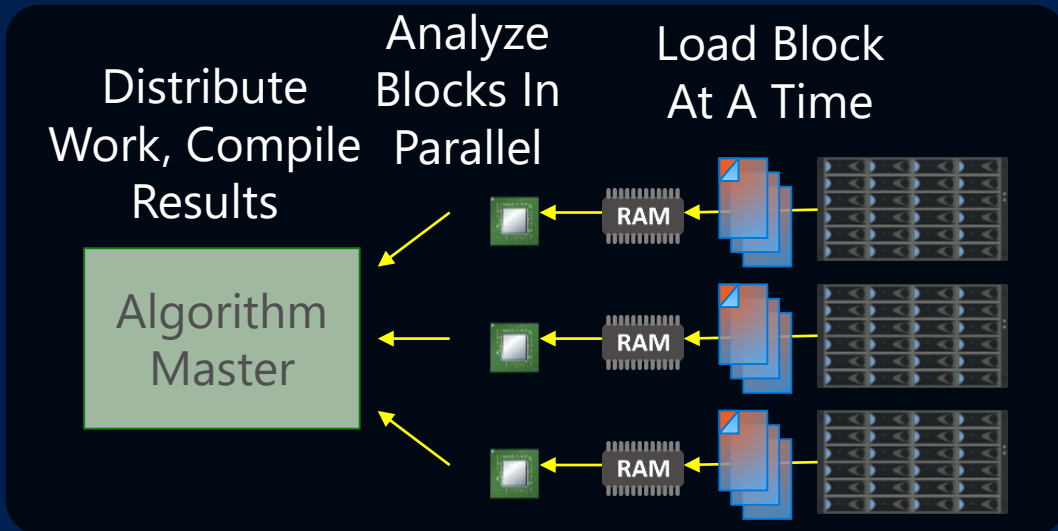
# Parallel Work Simple Example: Mean



Not every algorithm works in parallel  
Often there are several steps involved

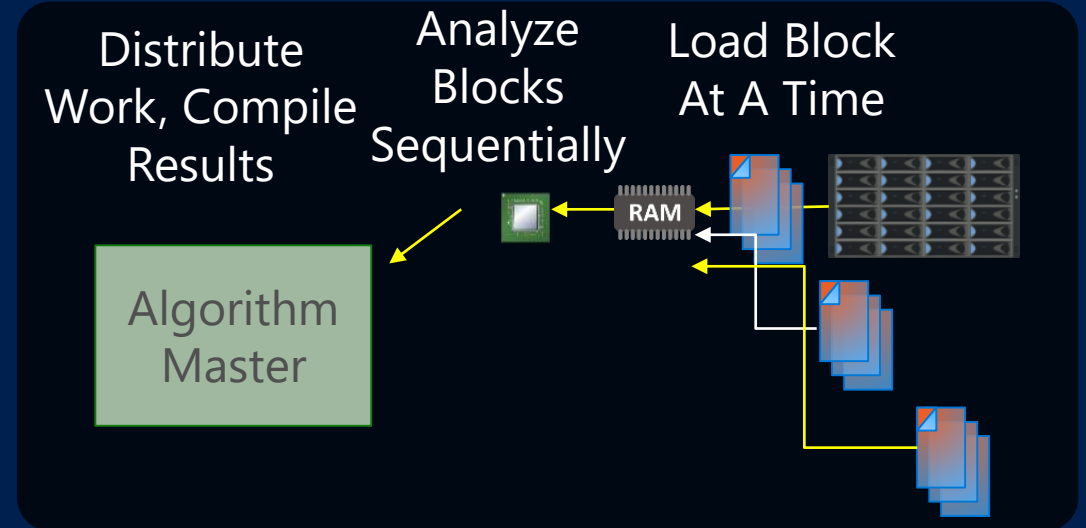


# Parallelization or Sequential Execution



## Parallel:

- Large number of parallel workers
- Fast handling of Big Data



## Sequential:

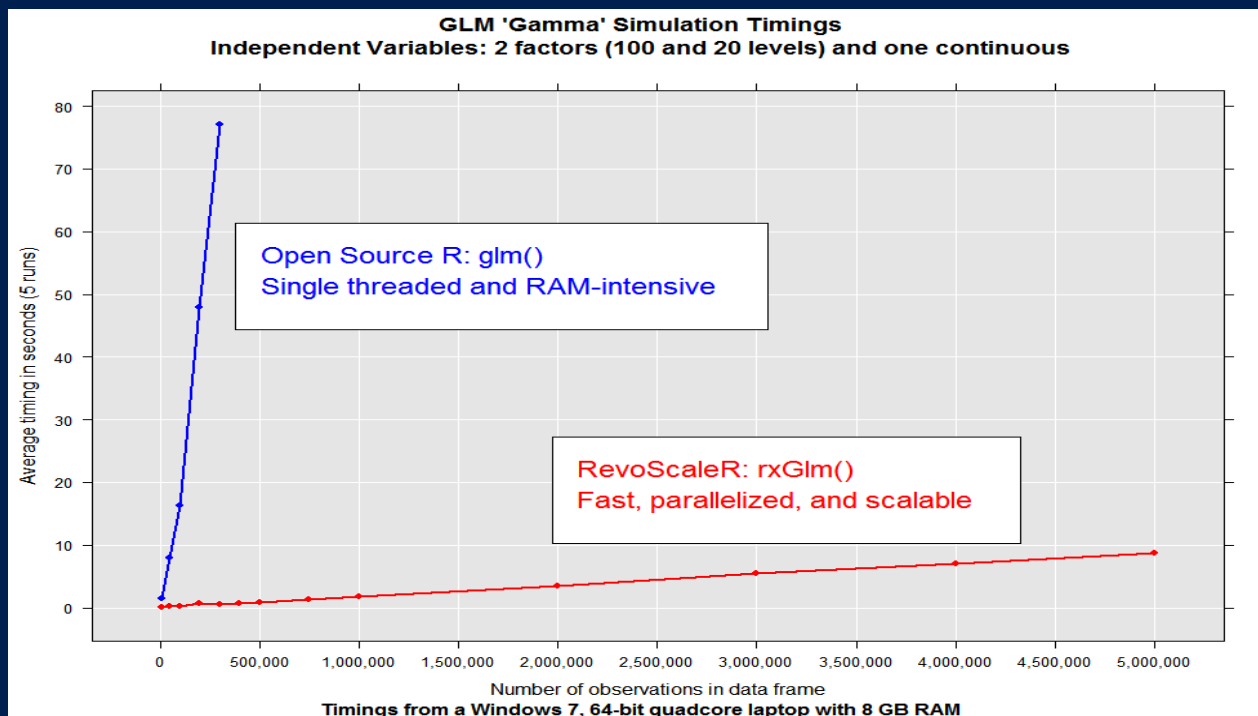
- Just one block of data in RAM
- Handling of Big Data with moderate resources



# Demo Microsoft R Server

# ScaleR - Performance comparison

Microsoft R Server has no data size limits in relation to size of available RAM

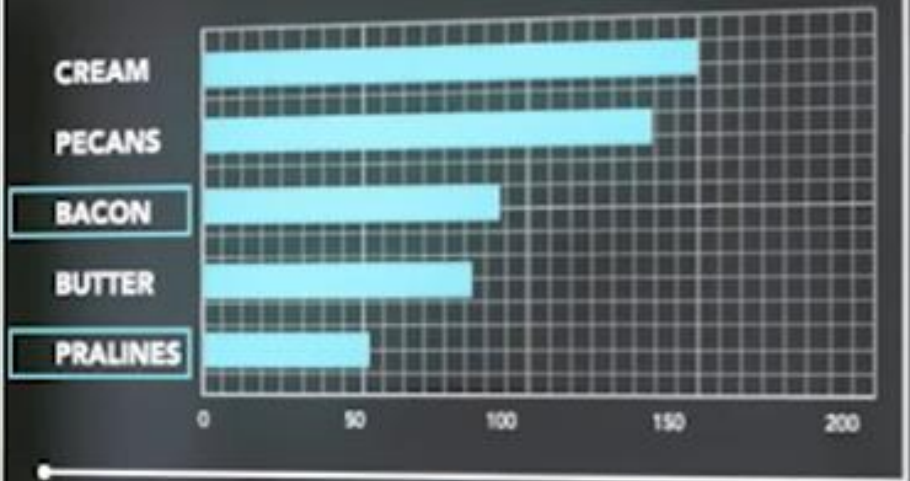


File Name	Compressed File Size (MB)	No. Rows	Open Source R (secs)	Revolution R (secs)
Tiny	0.3	1,235	0.00	0.05
V. Small	0.4	12,353	0.21	0.05
Small	1.3	123,534	0.03	0.03
Medium	10.7	1,235,349	1.94	0.08
Large	104.5	12,353,496	60.69	0.42
Big (full)	12,960.0	123,534,969	Memory!	4.89
V.Big	25,919.7	247,069,938	Memory!	9.49
Huge	51,840.2	494,139,876	Memory!	18.92

- US flight data for 20 years
- Linear Regression on Arrival Delay
- Run on 4 core laptop, 16GB RAM and 500GB SSD

**BEST SELLER:  
PECANS & CREAM**

□ SOCIAL AFFINITY SEARCH



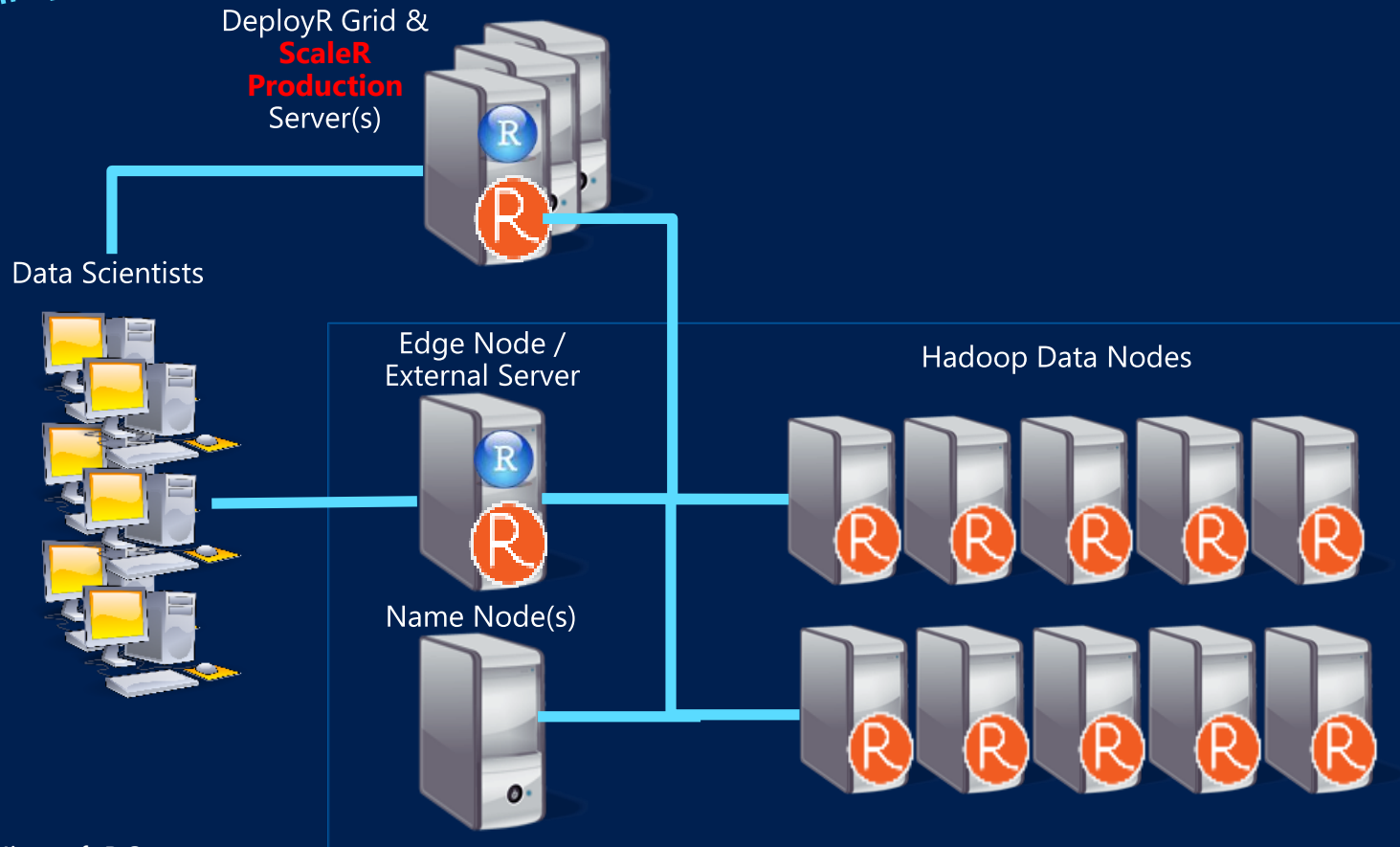
**SENTIMENT ANALYSIS: BACON + PRALINES**



Microsoft R Server  
Hadoop

# MRS and Hadoop Architecture options

## Conceptual Architecture



## Microsoft R model run options:

1. **Copy** HDFS data to Edge Node / External Server Linux file system. Use "**Local Parallel**" compute context on that server to run model
2. **Stream** data from HDFS to Edge Node / External Server. Use "**Local Parallel**" compute context on that server to run model and discard data
3. **Send** the Microsoft R model script to run in every data node and return model output object to the Edge Node / External Server

# Write Once Deploy Anywhere

ScaleR functions can run in-Hadoop, in-Spark or in-Database without any functional R recoding

Local Parallel – Linux or Windows

```
# SETUP LINUX ENVIRONMENT VARIABLES
rxSetComputeContext("localpar")

# CREATE LINUX, DIRECTORY AND FILE
OBJECTS
linuxFS <- RxNativeFileSystem()

AirlineDataSet <-
RxXdfData("AirlineDemoSmall.xdf",
fileSystem = linuxFS)
```

In – Hadoop

```
### SETUP HADOOP ENVIRONMENT VARIABLES
myHadoopCluster <- RxHadoopMR()

### HADOOP COMPUTE CONTEXT USING HDFS
rxSetComputeContext(myHadoopCluster)

### CREATE HDFS, DIRECTORY AND FILE OBJECTS
hdfsFS <- RxHdfsFileSystem()
AirlineDataSet <-
RxXdfData("AirlineDemoSmall.xdf",
fileSystem = hdfsFS)
```

And  
Spark

SQL Server

```
# SETUP MSSQLSERVER ENVIRONMENT VARIABLES
mySqlServer <- RxInSqlServer()

# SQL SERVER COMPUTE CONTEXT AND TABLE REF
rxSetComputeContext(mySqlServer)

AirlineDataSet <-
RxSqlServerData(table="AirlineDemoSmall")
```

And  
Teradata

R script – does not  
need to change to  
run across different  
platforms

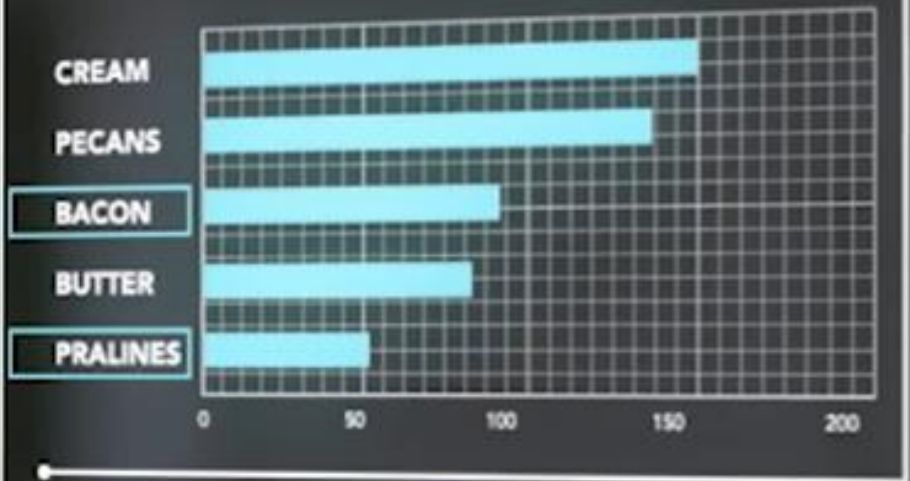
```
### ANALYTICAL PROCESSING ###
### Statistical Summary of the data
rxSummary(~ArrDelay+DayOfWeek, data= AirlineDataSet, reportProgress=1)

### CrossTab the data
rxCrossTabs(ArrDelay ~ DayOfWeek, data= AirlineDataSet, means=T)

### Linear Model and plot
hdfsXdfArrLateLinMod <- rxLinMod(ArrDelay ~ DayOfWeek + 0 , data =
AirlineDataSet)
plot(hdfsXdfArrLateLinMod$coefficients)
```

**BEST SELLER:  
PECANS & CREAM**

▣ SOCIAL AFFINITY SEARCH



**SENTIMENT ANALYSIS: BACON + PRALINES**



Microsoft R Server  
SQL Server R Services

# Why In-Database Analytics with SQL 2016 & R?

Leverage Full Capability of R:

- Rich Statistical, Visualization & Predictive Analytics
- A Large and Growing Skill Base

... including Microsoft R Servers Big Data Capabilities:

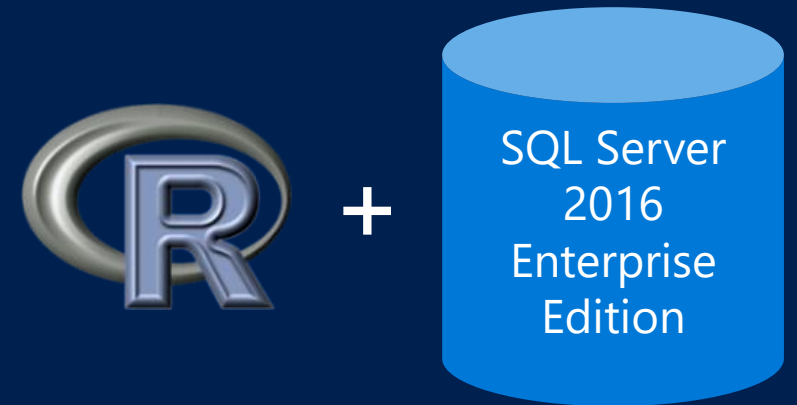
- Scalable Computation
- Scalable Data Size

... all Running In-Database:

- Divide Work Between Data Scientists and Data Engineers
- Reduce Data Duplication and Data Movement

... While Protecting Information:

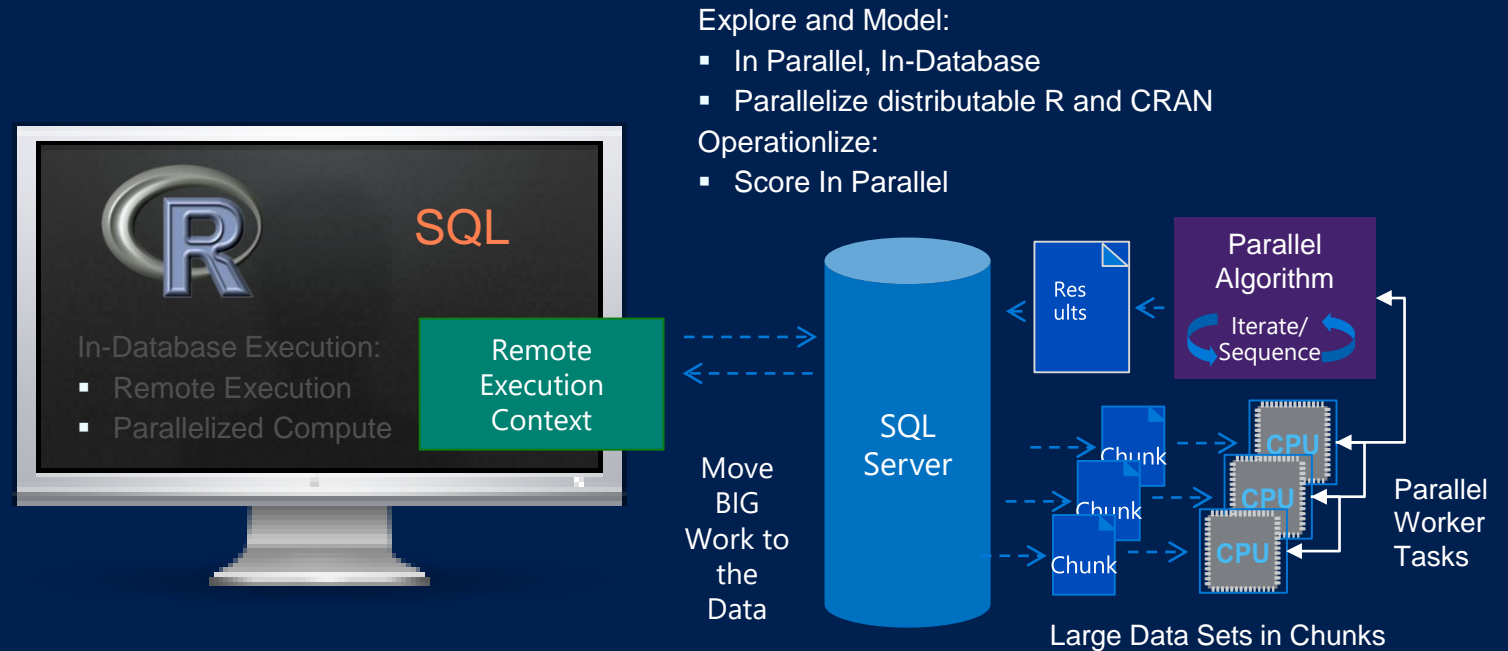
- Eliminate Data Movement & Unnecessary Copying
- Leverage Database Data Protections
- Leverage Database Tools for Backup, Scheduling, ...





# Run Parallel Algorithms in Database from an R client

- IDE for R
  - For Data Scientists
  - R coding
- SQL statements for data access
- SQL compute context
- Know-how:
  - R developer
  - Data science



## Explore and Model:

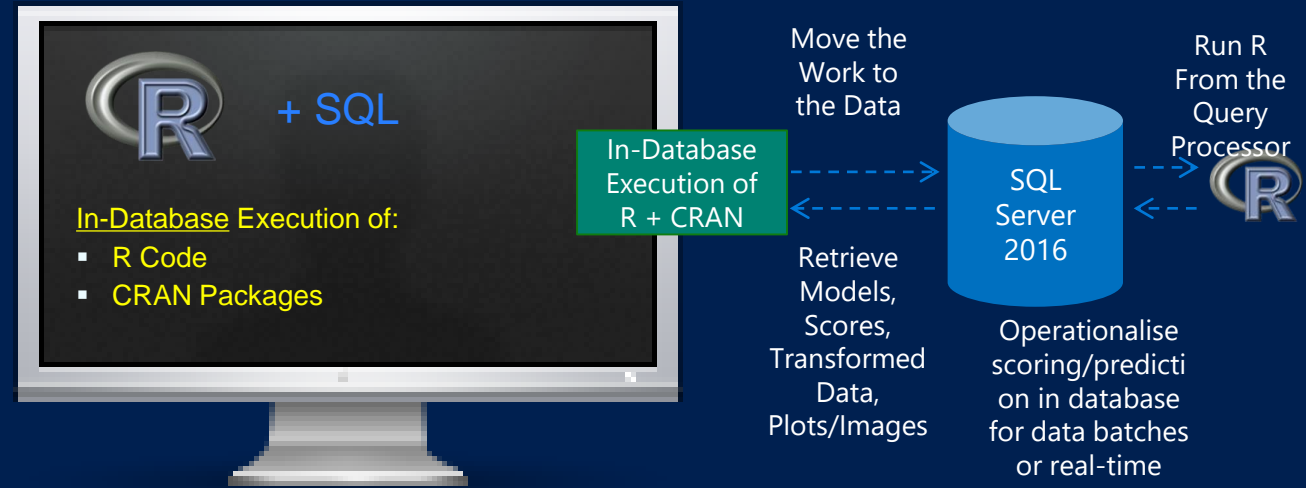
- In Parallel, In-Database
- Parallelize distributable R and CRAN

## Operationalize:

- Score In Parallel

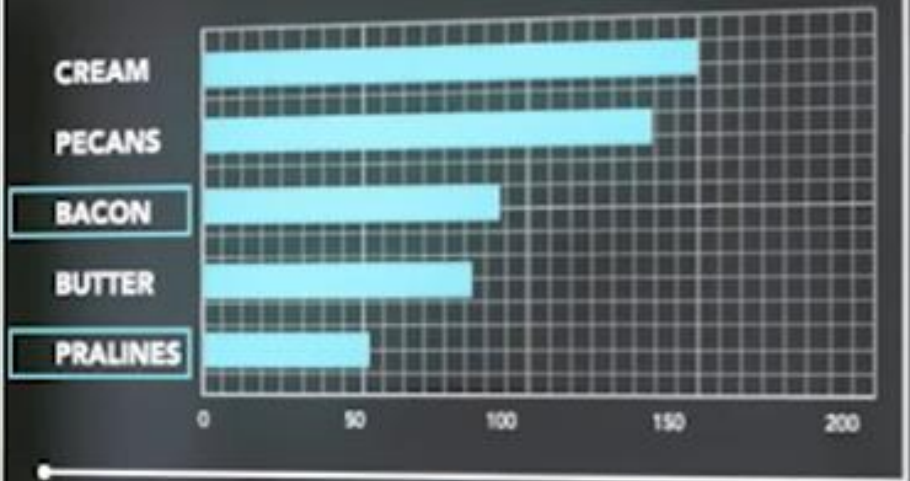
# Run R In-Database from TSQL

- IDE for SQL
  - For SQL developers, DB admins
- R code in stored procedures
- Know-how:
  - SQL developer
  - Operationalizing of SQL
  - Back-up, security, access control



**BEST SELLER:  
PECANS & CREAM**

□ SOCIAL AFFINITY SEARCH



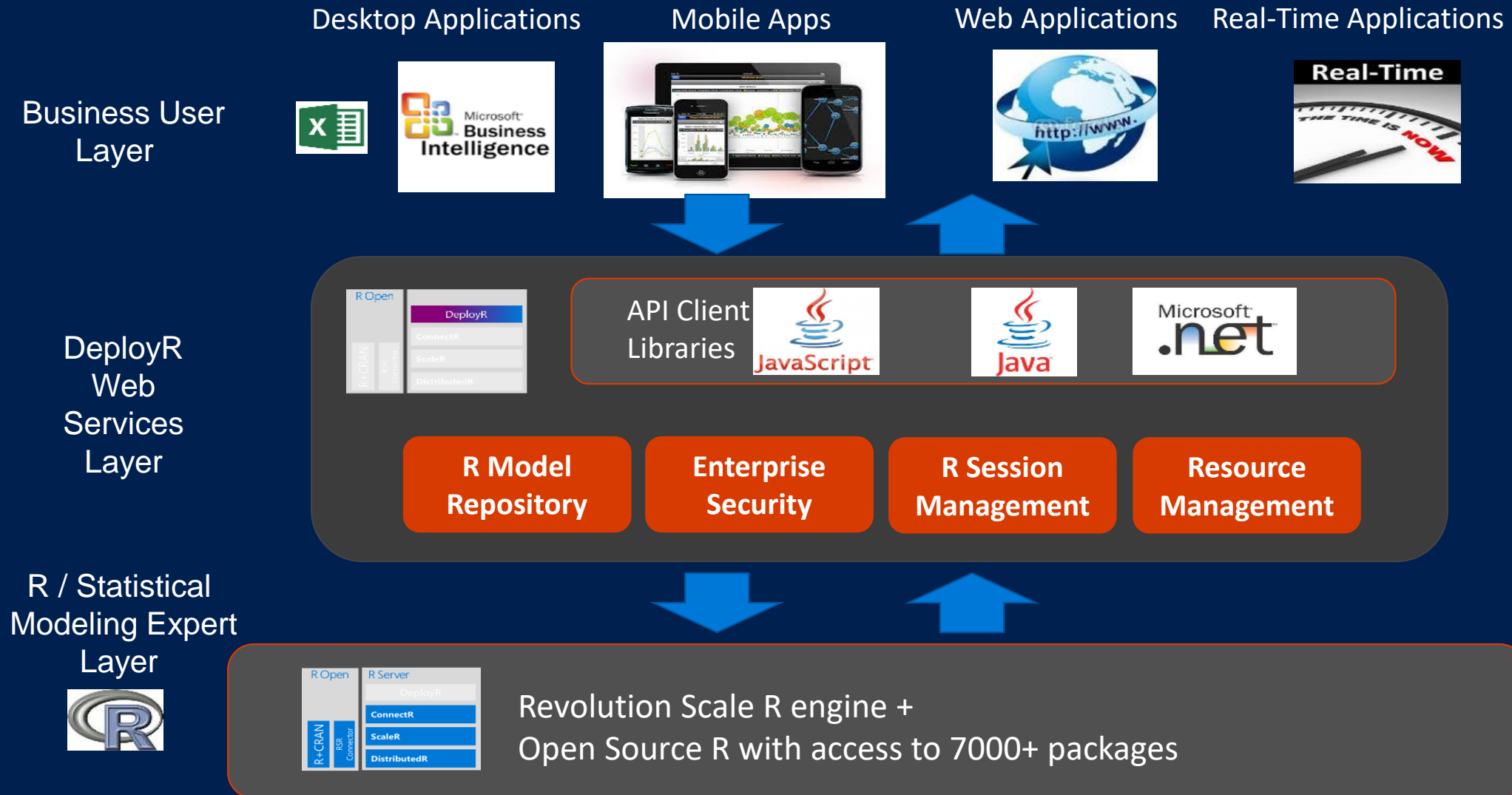
**SENTIMENT ANALYSIS: BACON + PRALINES**



Microsoft R Server

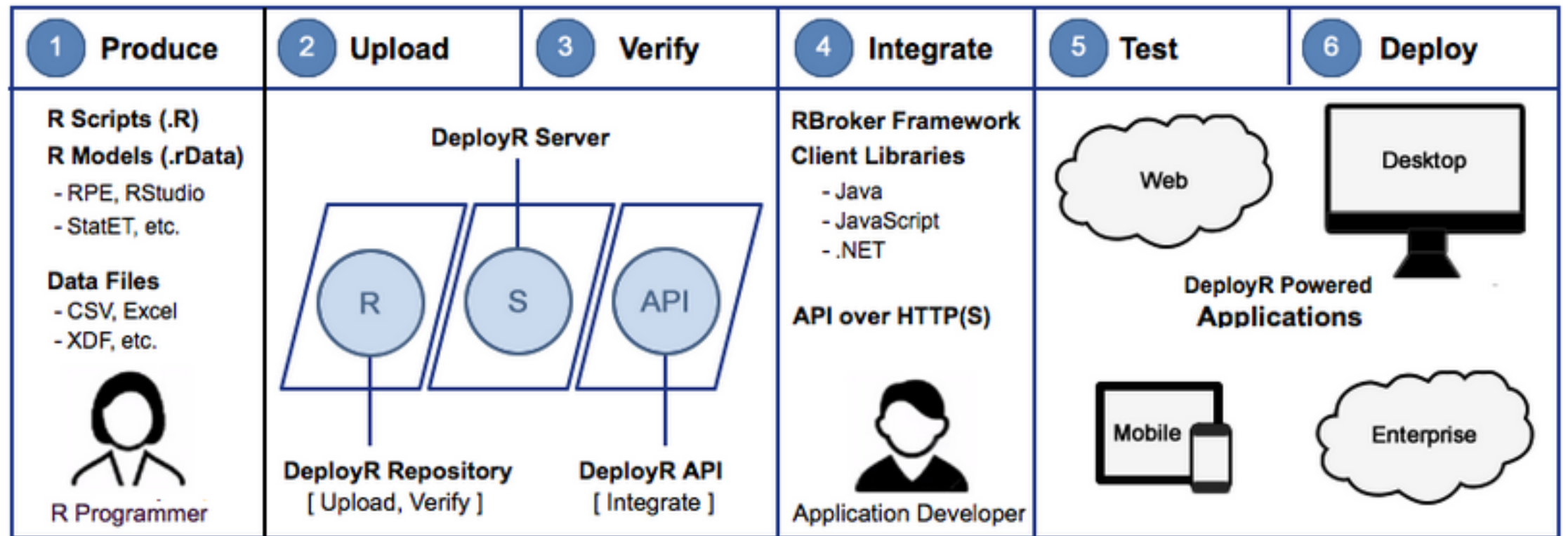
Operationalizing R based  
Analytics

# DeployR: Framework for R as a service for BI / web apps



# DeployR Workflow

## Basic Workflow



# DeployR – Administration

The creation and management of R runtime resource usage

**R Boundary List:** [+ New Boundary](#) [Export Boundaries](#) [Import Boundaries](#)

R boundaries constrain runtime resource usage related to user authenticated, asynchronous and anonymous operations on the grid.

R boundaries can be associated with any grid node or user. A default R boundary can be applied system-wide within Server Policies.

Grid node boundaries takes precedence over boundaries associated with users. If no other boundary is specified any default system-wide boundary is applied.

Name	Max CPU (sec)	Date Modified
DeployR Default R Boundary	1,740	1:19 PM on D
Boundary A	900	6:59 PM on M

The monitoring of events on the grid and server

[Home](#) [Users](#) [Roles](#) [R Scripts](#) [R Boundaries](#) [IP Filters](#) [The Grid](#) [Server Policies](#)

**Grid Node List:** [+ New Grid Node](#) [Export Grid Configuration](#) [Import Grid Configuration](#)

The Grid manages a scalable network of R processing nodes utilized by the RevoDeployR server. The processing resources associated with each node on the Grid can be dynamically allocated to specific modes of server operation including authenticated, asynchronous, anonymous and mixed-mode operations. This mechanism delivers a highly scalable, robust and secure framework for intensive R-compute environments.

Node Name	Description	Host	Operating Type	Slot Capacity	Enabled
DeployR Default Node	Default R processing node, mixed function.	localhost	Mixed Mode	Max: 100 @ 0.0%	true
Grid Node A	Authenticated Grid Node	162.209.10.235	Authenticated	Max: 100 @ 0.0%	true
Grid Node B	Asynchronous Grid Node	162.209.10.205	Asynchronous	Max: 100 @ 0.0%	true
Grid Node C	Anonymous Grid Node	162.209.10.208	Anonymous	Max: 100 @ 0.0%	true

**Grid Node List:** [+ New Grid Node](#) [Export Grid Configuration](#) [Import Grid Configuration](#)



# Azure Munich Meetup

Startseite Mitglieder Fotos Diskussionen Mehr Mein Profil



München, Deutschland  
Gegründet 23. Jul 2016

Über uns...

+ Freunde einladen

Mitglieder 161  
Anstehende 2

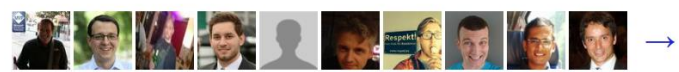
## Herzlich willkommen!

+ Schlage ein neues Meetup vor

Anstehende (2) Vergangene Kalender

### Azure IoT Development and Azure IoT Security

Microsoft Deutschland GmbH  
Walter-Gropius-Str. 5, München ([Karte](#))



TBD [weitere Infos](#)

Veranstaltet von: [Christian Waha](#) (Organisator)

Di, 22. Nov  
19:00

[RSVP](#)

43 gehen hin  
0 Kommentare

### IoT und Industrie 4.0: Von Buzzwords zu neuen Geschäftsmodellen

Microsoft Deutschland GmbH  
Walter-Gropius-Str. 5, München ([Karte](#))

Di, 13. Dez  
19:00

## Was gibt es Neues

NEUES MITGLIED  
[Martin Meixger](#)  
macht mit  
vor 2 Stunden

NEUES MITGLIED  
[Christian Straus](#)  
macht mit  
Gestern

NEUES RSVP  
[Marcus Franzen](#)  
hat Ja für IoT und Industrie 4.0: Von Buzzwords zu neuen Geschäftsmodellen geantwortet  
Vor 2 Tagen

NEUES RSVP  
[Marcus Franzen](#)  
hat Ja für Azure IoT Development and Azure



all azure communities?

<http://aka.ms/azure-meetups>