

Data Management Challenges in Quant Research: Solutions from OneTick and R

**Maria Belianina, Ph.D.
Director, Pre-Sales Engineering Support**

Contents

- ❑ **Data Management Challenges in Quant Research**
 - ❑ **R and OneTick: Addressing the challenges**
 - What is OneTick
 - Leveraging strong points of OneTick and R
 - OneTick data pre-processing and **R** analytics
 - **R** ↔ **OneTick** integration: 2 methods
 - ❑ **Examples**
 - Bring OneTick Queries or Raw Data into **R** environment
 - Use **R** Functions in OneTick queries
-

Data Management Challenges

Getting, storing and processing market data for quantitative research became more demanding as we face:

- ❑ **Increasing data granularity**
 - Daily to intraday
 - Milli → Micro → Nano → Picoseconds...
 - ❑ **Data cleansing** challenges
 - ❑ **Complexity** of data and data consolidation
 - Order book analytics
 - Order books vs Trades vs Custom and Summary data
 - ❑ Increasing data **volumes**
 - ❑ **Reference data** (corporate actions, name changes, continuous contracts, etc)
 - ❑ **Security master** maintenance
 - ❑ Database **schema changes**
-
- Reporting and analytics
vs all of the above**
-

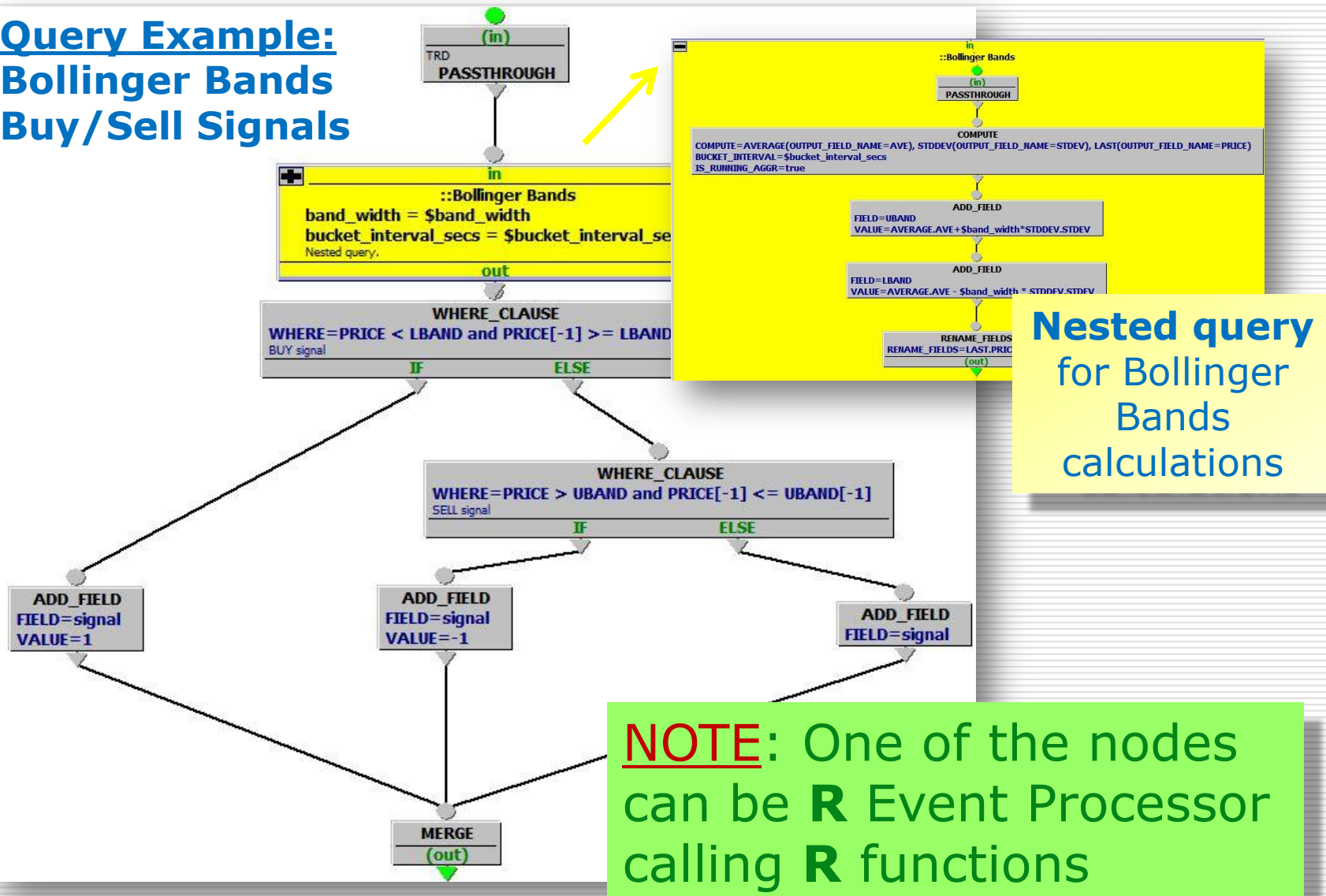
R AND ONETICK: ADDRESSING THE CHALLENGES

What is OneTick: Overview

- ✓ **Data feed adapters**, collectors and normalization
 - ✓ **Unlimited time series storage** for any security types and any tick types
 - ✓ In-memory database for **intraday ticks**
 - ✓ Complex Event Processing (**CEP**) engine for **real time analysis** & **signal generation**
 - ✓ **GUI** with sophisticated extendable **analytics** that can be integrated with **R** & **MatLab**
 - ✓ **API** (C++, C#, Java, Perl)
-

What is OneTick: GUI Analytics

Query Example: Bollinger Bands Buy/Sell Signals



What is OneTick: GUI Analytics

Query Example: Bollinger Bands Buy/Sell Signals



Symbol	Time	AVERAGE.AVE	STDDEV.STDEV	PRICE	UBAND	LBAND	signal	
O::MSFT.	2010/03/01 15:25:45.065	28.5800000000000	0.0000000000000	28.5800000000000	28.5800000000000	28.5800000000000	0	CI
O::MSFT.	2010/03/01 15:25:45.066	28.5800000000000	0.0000000000000	28.5800000000000	28.5800000000000	28.5800000000000	0	CI
O::MSFT.	2010/03/01 15:25:45.075	28.5800000000000	0.0000000000000	28.5800000000000	28.5800000000000	28.5800000000000	0	CI
O::MSFT.	2010/03/01 15:25:45.117	28.5800000000000	0.0000000000000	28.5800000000000	28.5800000000000	28.5800000000000	0	CI
O::MSFT.	2010/03/01 15:25:45.257	28.5840000000000	0.0080000000000	28.6000000000000	28.5840800000000	28.5839200000000	-1	CI
O::MSFT.	2010/03/01 15:25:45.281	28.5866666666666	0.009428090415	28.6000000000000	28.586760947570	28.586572385762	0	CI

NOTE: Query can be called from **R** passing this output to **R** vector

OneTick and Market Data Management

Market Data Challenges...	OneTick Solution For All Security Types Include...
<ul style="list-style-type: none"> Increasing volumes History vs Real Time 	<ul style="list-style-type: none"> Archives are <u>compressed</u>, <u>distributed</u>, <u>scalable</u> <u>Intraday in-memory DB</u>+ <u>CEP Real-time</u> + <u>Archives</u>
<ul style="list-style-type: none"> Increasing granularity Changing schemas Reference data Master data 	<ul style="list-style-type: none"> Time series: <u>flexible schema</u> & <u>nanoseconds</u> <u>Optimized summary DBs</u> for low frequency data <u>Reference data</u> design/analytics (naming continuity, contract rollovers, calendars, corporate actions)
<ul style="list-style-type: none"> Data cleansing 	<ul style="list-style-type: none"> Normalization rules start the <u>cleansing process</u> Data can be <u>updated</u> or <u>deleted</u> by queries Parts of archives & in-memory DB can be updated
<ul style="list-style-type: none"> Complexity Consolidation Order books 	<ul style="list-style-type: none"> Built-in <u>order book storage and analytics</u> Order book data can be <u>joined</u> or <u>merged</u> with any other kind of data within analytical queries Join data for derivatives and underlyers
<ul style="list-style-type: none"> Reporting & analytics vs all of the above 	<ul style="list-style-type: none"> <u>Rich GUI and API analytics</u> <u>DB structure is transparent</u> to queries Queries for <u>1,000s of symbols</u> & <u>multiple sources</u> <u>Integration</u>: R, MatLab, SAS, C++, Java, Perl

System Integration: OneTick and R

	OneTick allows to:	R allows to:
If you work in R:	Pre-process data in OT: <ul style="list-style-type: none">✓ <u>Normalize</u> and clean data✓ Apply <u>reference data</u>✓ <u>Aggregate</u> or <u>intervalize</u> irregular time series✓ <u>JOIN</u> or <u>MERGE</u> trades with order books, news, etc...✓ <u>Use OT functions</u> to preprocess✓ <u>Parameterize</u> queries	<ul style="list-style-type: none">✓ Retrieve OT query output in R code via ODBC/OneTick SQL:<ul style="list-style-type: none">• <u>Call tested OT query</u>• <u>Pass parameters</u>• Results go into an R <u>vector</u>• Process results vector✓ <u>Limit amount of data</u> processed within ODBC and R
If you work in OneTick GUI:	<ul style="list-style-type: none">✓ [All of the above]✓ Call any number of OT/R functions✓ Call <u>separate R functions</u> for the same tick bucket for bucket <u>entry</u> or <u>exit events</u>✓ Run <u>historical</u> or <u>CEP queries</u>✓ <u>View query results</u> in GUI grid or chart, or via API	<ul style="list-style-type: none">✓ <u>Map OneTick fields</u> to R parameters✓ <u>Map R output</u> back to OneTick fields for further OneTick query processing <hr/> Tip: Take full advantage of both packages.

Integration Method 1:

**ONETICK QUERY RESULTS →
BACK TO R**

OneTick Results Back to R: Overview

Prerequisites: OneTick client setup

Connection: OneTick ODBC Driver

Syntaxes: OneTick SQL (based on SQL)

OneTick Query prep steps:

- ✓ Design query graph with OneTick GUI
- ✓ Create query parameters to be passed from **R**
- ✓ Test and save the query

Results: **R** vector

Via OT ODBC + SQL call with parameters

OneTick Results Back to R: Overview

- P** **TIP:**
C **Pack logic and analytics**
S **into OneTick query**
O to take advantage of OneTick Server,
and
to limit network traffic and
a number of returned ticks
- ✓ Create query parameters to be passed from R
 - ✓ Test and save the query
- Results: R vector**
Via OT ODBC + SQL call with parameters
-

OneTick Results Back to R: Example

Sample Business Case:

I. OneTick Query must produce:

Running liquidity profile (average daily volume over X past days) **grouped by "Time of Day"**.

E.g., average 5-day volume for 9:30 – 10:30, timestamp'ed by the 5th day

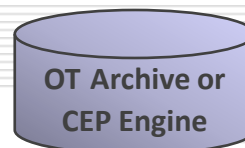
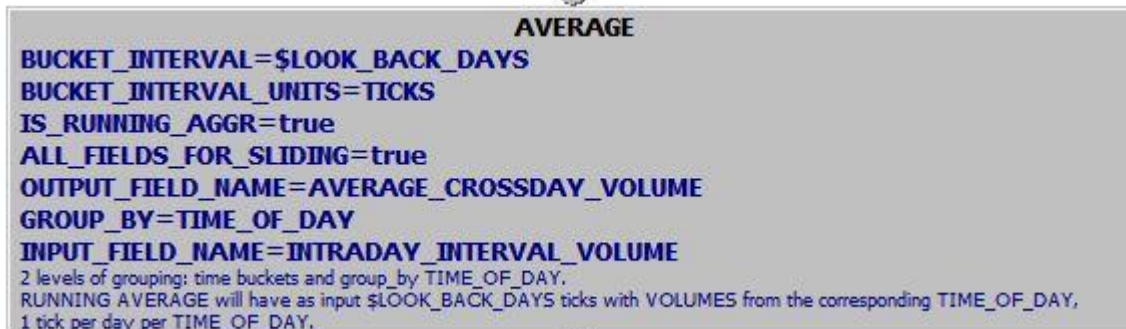
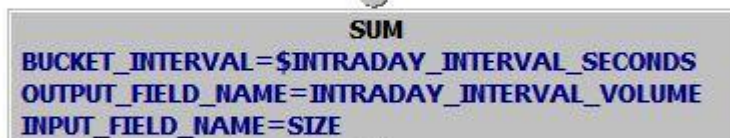
Source trade data: 1 month, 1 symbol

II. Results → Back to R

for plotting and further processing

OneTick Results Back to R: Design OT Query

OT GUI Query Graph Design



Tick Processing

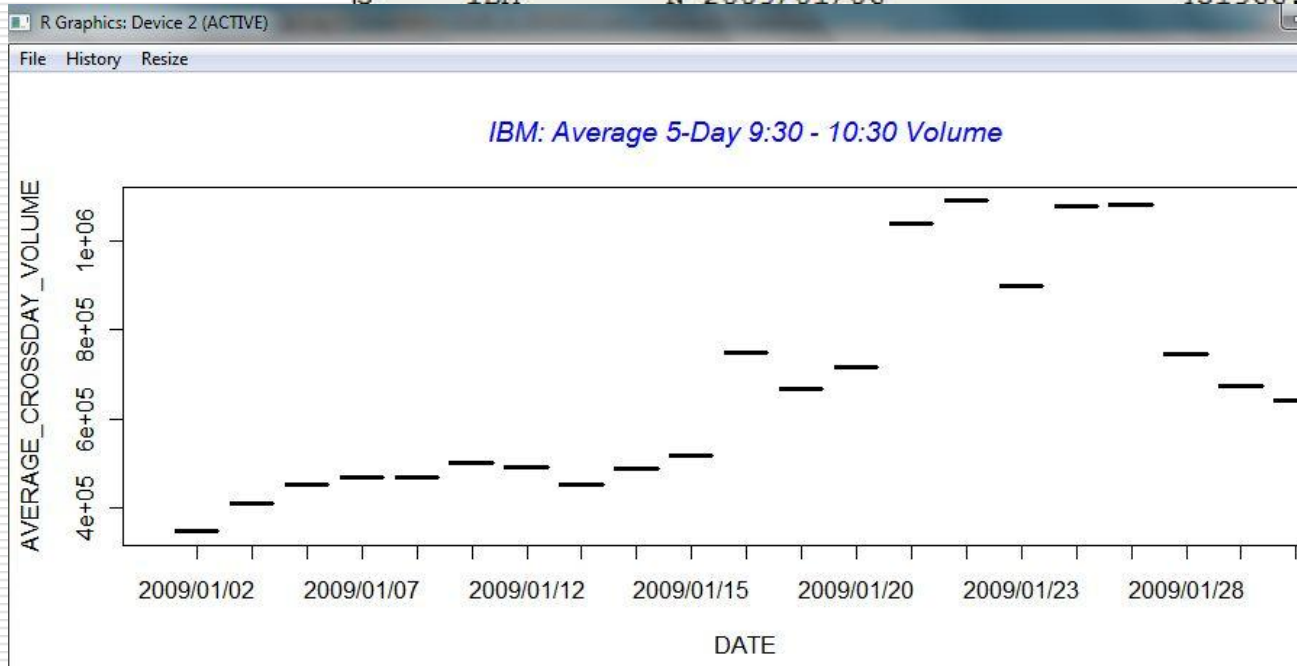
Approach & OT Steps

- Retrieve trades, filter
- Calculate **VOLUMEs** for all INTRADAY INTERVALS (in seconds)
- Add TIME_OF_DAY fields to be used for grouping
- Calculate **running AVERAGE**
 - for each TIME_OF_DAY
 - over LOOK_BACK_DAYS for **LIQUIDITY PROFILE** for each day in query date range
- **Test** query in GUI

OneTick Results Back to R: Results 1

- ✓ Pass SQL
- ✓ Get query results for 1 month, 9:30 - 10:30 only
- ✓ Plot 5-day volume average by day

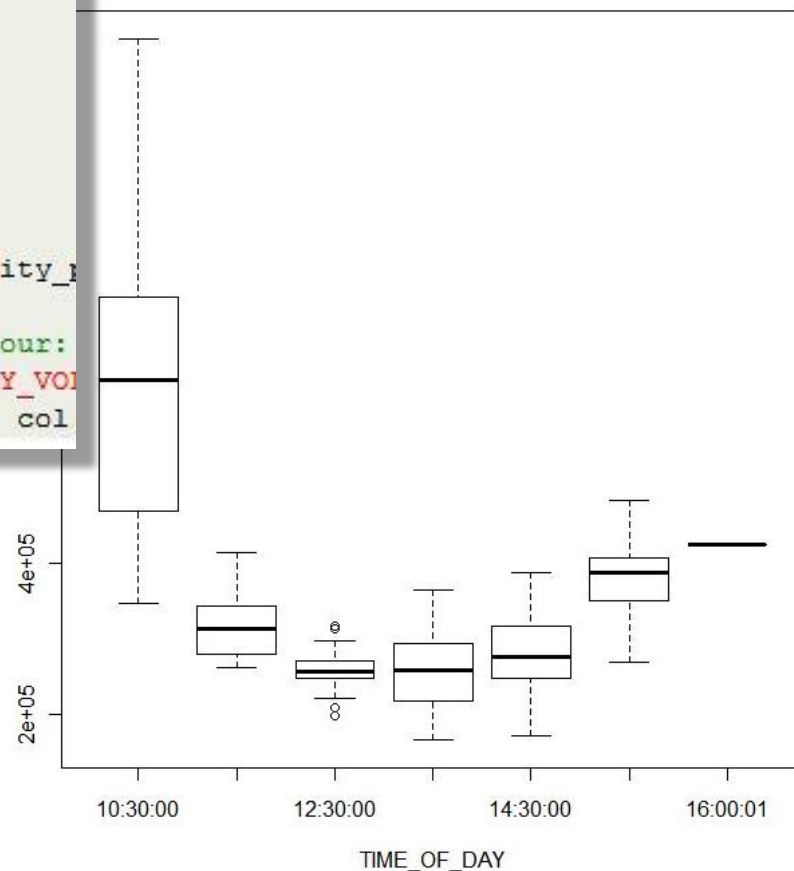
```
REvolution R Enterprise Console
> # Print out selected columns:
> results[, c("SYMBOL", "EXCHANGE", "DATE", "AVERAGE_CROSSDAY_VOLUME")]
  SYMBOL EXCHANGE      DATE AVERAGE_CROSSDAY_VOLUME
1    IBM         N 2009/01/02             346800.0
2    IBM         N 2009/01/05             411000.0
3    IBM         N 2009/01/06            451966.7
```



OneTick Results Back to R: Results 2

```
REvolution R Enterprise Console
> # Create a new chart:
> sql<-omdBuildSQLQuery("sample_liquidity_profile_R",
+ "2009-01-02 09:30:00",
+ "2009-01-30 16:00:01",
+ "EST5EDT",
+ DB="TAQ",
+ SYMBOL="IBM.N",
+ INTRADAY_INTERVAL_SECONDS=3600,|
+ LOOK_BACK_DAYS=5,
+ TIME_OF_DAY='ALL'
+ ); sql
[1] "SELECT * FROM OMD.OTQ_FILES.\"sample_liquidity_p
> results<-sqlQuery(channel, sql);
> # View behavior of crossday volume average by hour:
> plot(results[, c("TIME_OF_DAY", "AVERAGE_CROSSDAY_VOI
> title(main="IBM: Average 5-Day Volume by Hour", col
```

IBM: Average 5-Day Volume by Hour



- ✓ Pass SQL
- ✓ Get query results for 1 month
- ✓ Plot 5-day volume average values by TIME_OF_DAY

Integration Method 2:

ONETICK QUERY ↔
R FUNCTIONS

R Functions in OneTick Queries: Overview

Prerequisites: R or REvolution R installation on the OneTick server (depends on OS)

Connection: Using standard R DLL

Syntaxes: OneTick GUI and R expressions

Input: OneTick Archive or real-time tick timeseries from a single or multiple data sources

Output: OneTick query results
(timeseries defined by the query design)

Query Types: Historical or Continuous CEP

R Functions in OneTick: Example 1

Business Case:

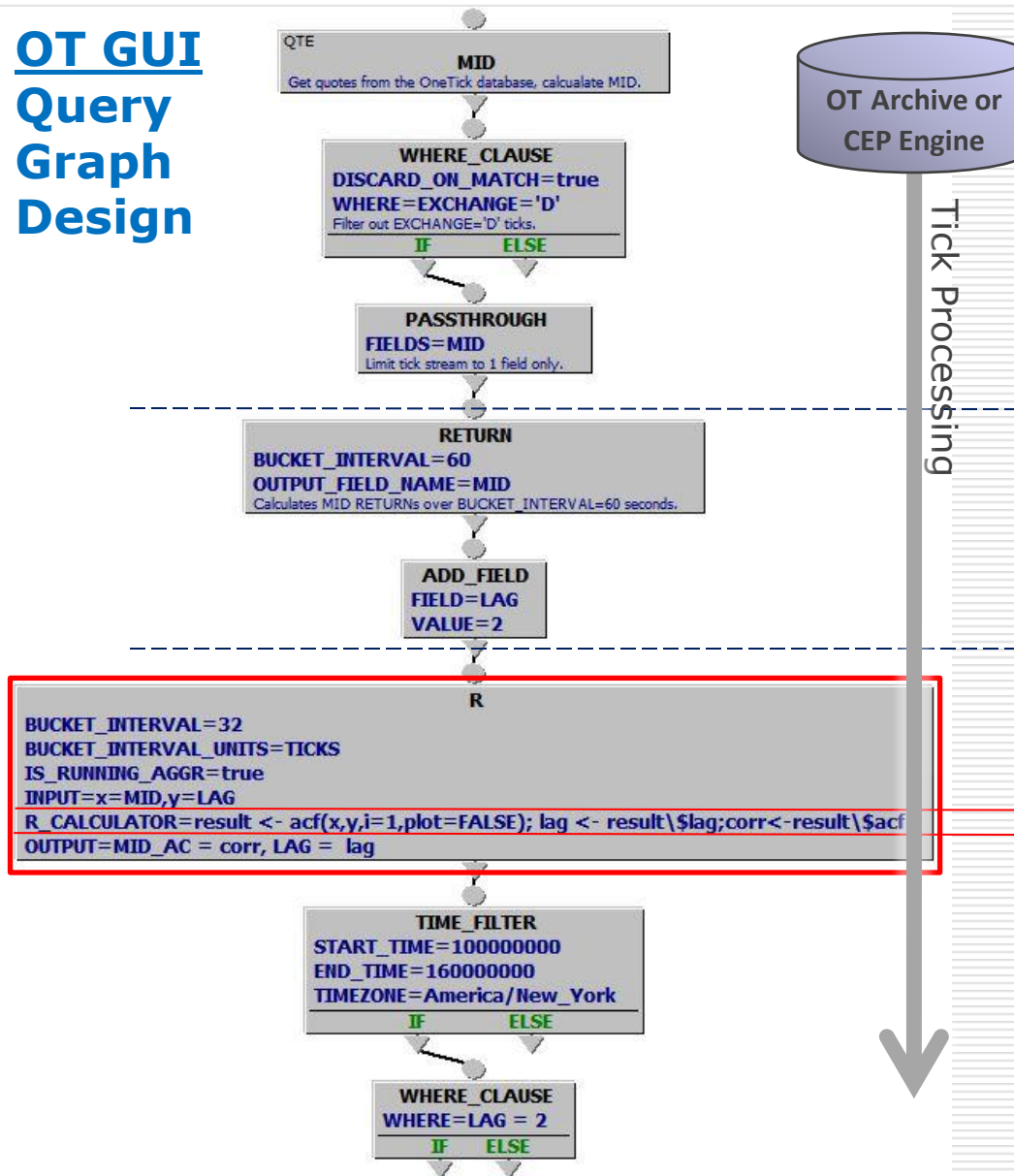
Find serial **autocorrelation**
of 30 one-minute returns with lag of 2

R Function:

acf, an autocorrelation function

R Functions in OneTick: Design Query

OT GUI Query Graph Design



Approach & OT Steps

- Retrieve quote ticks for any number of symbols
- Add tick by tick MID calculation
- Filter ticks

- Aggregate ticks every 60 seconds
- Calculate **RETURN(MID)**
- Add LAG field to each tick

- Create running (a.k.a. sliding) aggregation of 32 ticks
- Call **R** function **acf(...)** for each sliding group
- Pass values of MID and LAG

- Filter resulting ticks
- Review results

R Functions in OneTick: R Processor

OT GUI
Query
Graph
Design

R parameters

BUCKET_INTERVAL	32	secs/ticks
BUCKET_INTERVAL_UNITS	TICKS	
OUTPUT_INTERVAL		seconds
OUTPUT_INTERVAL_UNITS	SECONDS	
IS_RUNNING_AGGR	true	
BUCKET_TIME	BUCKET_END	
BUCKET_END_CRITERIA		
ALL_FIELDS_FOR_SLIDING	false	
PARTIAL_BUCKET_HANDLING	AS_SEPARATE_BUCKET	

R out-of-box OneTick event processor
based on **R** distributed dll.

Can be used multiple times in 1 query with different R functions

R_INITIALIZER		Expression in R
INPUT	x=MID,y=LAG	Mapping from tick
R_CALCULATOR	result <- acf(x,y,i=1,plot=FALSE); lag <- result\$lag; corr <- result\$acf	Expression in R
OUTPUT	MID_AC = corr, LAG = lag	Mapping from R
R_ENTERING_TICK_HANDLER		For running
R_LEAVING_TICK_HANDLER		For running

R function call parameters

GROUP_BY

Description Node name Tick type

TIMEZONE=America/New_York

IF ELSE

WHERE_CLAUSE
WHERE=LAG = 2
IF ELSE

Filter resulting ticks

- Review results

BUCKET_INTERVAL=32
BUCKET_INTERVAL_UNITS=
IS_RUNNING_AGGR=true
INPUT=x=MID,y=LAG
R_CALCULATOR=result <-
OUTPUT=MID_AC = corr, L

R Functions in OneTick: Query Output

View OT Query results in a GUI grid or chart

Index	Symbol	Time	LAG	MID_AC
1	FULL_DEMO_L1::CSCO	2006/06/01 10:00:00.000	2.00000000000000	-0.254566084117
2	FULL_DEMO_L1::CSCO	2006/06/01 10:01:00.000	2.00000000000000	-0.246575301054
3	FULL_DEMO_L1::CSCO	2006/06/01 10:02:00.000	2.00000000000000	-0.2427
4	FULL_DEMO_L1::CSCO	2006/06/01 10:03:00.000	2.00000000000000	-0.2470
5	FULL_DEMO_L1::CSCO	2006/06/01 10:04:00.000	2.00000000000000	-0.2540
6	FULL_DEMO_L1::CSCO	2006/06/01 10:05:00.000	2.00000000000000	-0.2534
7	FULL_DEMO_L1::CSCO	2006/06/01 10:06:00.000	2.00000000000000	-0.2434
8	FULL_DEMO_L1::CSCO	2006/06/01 10:07:00.000	2.00000000000000	-0.2362
9	FULL_DEMO_L1::CSCO	2006/06/01 10:08:00.000	2.00000000000000	-0.2339

Or pass query results
into API wrapper
application :
C++, C#, Java, Perl, or
a command line query

Date and time Start End Timezone

MM/DD/YYYY 06/01/2006 06/01/2006 New York Switch Datasets

HH:MM:SS.mmm 09:30:00.000 16:30:00.000 ☐ Apply times daily ☐ Continuous query

Start expression



R Functions in OneTick: R Parameters

2 subsets of parameters that work together:

OneTick aggregation	R function specifications
BUCKET_INTERVAL , UNITS and optional GROUP_BY to aggregate ticks into buckets	R_INITIALIZER to specify one time only initial command
OUTPUT_INTERVAL and UNITS to define frequency of output for running calculations	INPUT , OUTPUT , R_CALCULATOR to map tick fields to R variables and specify R command
IS_RUNNING = true/false for running (<i>a.k.a.</i> sliding) calculations Other parameters for additional flexibility with aggregation	R_ENTERING/LEAVING_TICK_HANDLER to call different functions for BUCKET_INTERVAL entry and exit events in RUNNING aggregations

Result:

R functions in **OneTick** queries at **any level of granularity**, with **OneTick** processing and aggregating **large datasets**, for **any number of symbols**, in **historical** or **real time** queries.

Q&A

Contact:

Dr. Maria Belianina

maria.belianina@onetick.com

support@onetick.com
