### Advanced Data Management (CSCI 640/490)

### Time Series Data

Dr. David Koop





### Dataframes, Databases, and the Cloud

- How do we take advantage of different architectures? Lots of work in scaling databases and specialized computational engines • What is the code that people actually write?

### D. Koop, CSCI 640/490, Spring 2023





2

### Data Science Jungle



D. Koop, CSCI 640/490, Spring 2023



Northern Illinois University



## Magpie Goals



### D. Koop, CSCI 640/490, Spring 2023



Northern Illinois University

NIU

### Magpie Architecture







### ConnectorX: Databases to Dataframes











### Dataframe API?

### • SQL, pandas, or something else?



**Doris Lee** @dorisilee

🥒 Hot Takes on Enterprise Pandas — Day 2 🌙



In many cases, SQL isn't the solution, and pandas is the easier path. below

11:15 AM · Mar 27, 2023 · 6,517 Views



**Doris Lee** @dorisjlee

SQL is good for certain things, but there are things that SQL wasn't meant to do, and if you contort SQL to do them, you wind up with nightmarish queries. Many of these can be no more than a few lines in pandas.

11:15 AM · Mar 27, 2023 · **172** Views

### D. Koop, CSCI 640/490, Spring 2023



🌙 Hot Takes on Enterprise Pandas — Day 5 🌙

Beware of "pandas-like" APIs that aren't actually compatible with pandas. Many dataframe libraries may look similar to pandas but lack support for critical pandas functionalities.

...

...

11:15 AM · Mar 30, 2023 · 5,225 Views





...



### <u>Assignment 4</u>

- Work on Data Integration and Data Fusion
- World Bank, OECD)
  - Integrate information with population
- Record Matching:
  - Which countries are the same?
- Data Fusion:
  - The receipts/expenditures
  - Country names

### D. Koop, CSCI 640/490, Spring 2023

# Integrate travel datasets from different institutions (UN World Tourism Office,







## Test 2

- Upcoming... April 10
- research papers

### • Similar format, but more emphasis on topics we have covered including the







### Time Series Data





### What is time series data?

- Technically, it's normal tabular data with a timestamp attached
- This allows more analysis
- Example: Web site database that tracks the last time a user logged in

  - 2: Add a new row with login information every time the user logs in
  - Option 2 takes more storage, but we can also do a lot more analysis!

But... we have observations of the same values over time, usually in order

- 1: Keep an attribute lastLogin that is overwritten every time user logs in





## What is Time Series Data?

• A row of data that consists of a timestamp, a value, optional tags

## timestamp

un						
time	generated	message_subtype	scaler	short_id	tenant	value
2016-07-12T11:51:45Z	"true"	"34"	"4"	"3"	"saarlouis"	465110000
2016-07-12T11:51:45Z	"true"	"34"	"-6"	"2"	"saarlouis"	0.061966999999999994
2016-07-12T12:10:00Z	"true"	"34"	"7"	"5"	"saarlouis"	49370000000
2016-07-12T12:10:00Z	"true"	"34"	"6"	"2"	"saarlouis"	18573000000
2016-07-12T12:10:00Z	"true"	"34"	"5"	"7"	"saarlouis"	5902300000

### D. Koop, CSCI 640/490, Spring 2023

### tags

value







12



### Time Series Data

- Metrics: measurements at regular intervals
- Events: measurements that are not gathered at regular intervals











### Types of Time Series Data

- time series: observations for a **single** entity at **different** time intervals - one patient's heart rate every minute
- cross-section: observations for **multiple** entities at the **same** point in time
- heart rates of 100 patients at 8:01pm
- panel data: observations for **multiple** entities at **different** time intervals - heart rates of 100 patients every minute over the past hour









### Features of Time Series Data

- Trend: long-term increase or decrease in the data
- Seasonal Pattern: time series is affected by seasonal factors such as the time of the year or the day of the week (fixed and of known frequency)
- Cyclic Pattern: rises and falls that are not of a fixed frequency
- Stationary: no predictable patterns (roughly horizontal with constant variance)
  - White noise series is stationary
  - Will look the basically the same whenever you observe it















































































## Types of Time Data

- Timestamps: specific instants in time (e.g. 2018 11 27 14:15:00)
- Periods: have a standard start and length (e.g. the month November 2018)
- Intervals: have a start and end timestamp
  - Periods are special case
  - Example: 2018-11-21 14:15:00 2018-12-01 05:15:00
- Elapsed time: measure of time relative to a start time (15 minutes)





### Dates and Times

- What is time to a computer?
  - Can be stored as seconds since Unix Epoch (January 1st, 1970)
- Often useful to break down into minutes, hours, days, months, years...
- Lots of different ways to write time:
  - How could you write "November 29, 2016"?
  - European vs. American ordering...
- What about time zones?





## Python Support for Time

- The datetime package
  - Has date, time, and datetime classes
  - .now() method: the current datetime
- Can access properties of the time (year, month, seconds, etc.) Converting from strings to datetimes:
  - datetime.strptime: good for known formats
  - dateutil.parser.parse: good for unknown formats
- Converting to strings
  - str(dt) Or dt.strftime(<format>)





### Datetime format specification

- Look it up:
  - <u>http://strftime.org</u>
- Generally, can create whatever format you need using these format strings

Code	Meaning	Example
%a	Weekday as locale's abbreviated name.	Mon
۶A	Weekday as locale's full name.	Monday
88	Weekday as a decimal number, where 0 is Sunday and 6 is Saturday.	1
%d	Day of the month as a zero-padded decimal number.	30
%−d	Day of the month as a decimal number. (Platform specific)	30
%b	Month as locale's abbreviated name.	Sep
۶B	Month as locale's full name.	September
%m	Month as a zero-padded decimal number.	09
%-m	Month as a decimal number. (Platform specific)	9
۶y	Year without century as a zero-padded decimal number.	13
8Y	Year with century as a decimal number.	2013
%H	Hour (24-hour clock) as a zero-padded decimal number.	07
%−H	Hour (24-hour clock) as a decimal number. (Platform specific)	7
%I	Hour (12-hour clock) as a zero-padded decimal number.	07
%-I	Hour (12-hour clock) as a decimal number. (Platform specific)	7
%p	Locale's equivalent of either AM or PM.	AM
۶M	Minute as a zero-padded decimal number.	06
%-M	Minute as a decimal number. (Platform specific)	6
°S	Second as a zero-padded decimal number.	05
%-S	Second as a decimal number. (Platform specific)	5









### Pandas Support for Datetime

- pd.to datetime:
  - convenience method
  - can convert an entire column to datetime
- Has a Nat to indicate a missing time value
- Stores in a numpy.datetime64 format
- pd.Timestamp: a wrapper for the datetime 64 objects









## More Pandas Support

- can be interpreted as a date:
  - ts['1/10/2011'] Or ts['20110110']
- Date ranges: pd.date range('4/1/2012', '6/1/2012', freq='4h')
- Slicing works as expected
- Can do operations (add, subtract) on data indexed by datetime and the indexes will match up
- As with strings, to treat a column as datetime, you can use the .dt accessor

### D. Koop, CSCI 640/490, Spring 2023

Accessing a particular time or checking equivalence allows any string that









## Generating Date Ranges

- index = pd.date range('4/1/2012', '6/1/2012')
- Can generate based on a number of periods as well - index = pd.date range('4/1/2012', periods=20)
- Frequency (freq) controls how the range is divided
  - Codes for specifying this (e.g. 4h, D, M)
  - In [90]: pd.date range('1/1/2000', '1/3/2000 23:59', freq='4h') Out[90]: <class 'pandas.tseries.index.DatetimeIndex'>  $[2000-01-01 \ 00:00:00, \ldots, 2000-01-03 \ 20:00:00]$ Length: 18, Freq: 4H, Timezone: None
  - Can also mix them: '2h30m'











## Time Series Frequencies

Alias	Offset Type
D	Day
В	BusinessDay
Н	Hour
T or min	Minute
S	Second
L or ms	Milli
U	Micro
Μ	MonthEnd
BM	BusinessMonthEnd
MS	MonthBegin
BMS	BusinessMonthBegin
W-MON, W-TUE,	Week
WOM-1MON, WOM-2MON,	WeekOfMonth

	Description	
	Calendar daily	
	Business daily	
	Hourly	
	Minutely	
	Secondly	
	Millisecond (1/1000th of 1 second)	
	Microsecond (1/1000000th of 1 second	)
	Last calendar day of month	
	Last business day (weekday) of month	
	First calendar day of month	
1	First weekday of month	
	Weekly on given day of week: MON, TU or SUN.	IE, WED, THU, FRI, SAT,
	Generate weekly dates in the first, secor of the month. For example, WOM-3FR	nd, third, or fourth week E for the 3rd Friday of
	each month.	[W. McKinney, F









### DatetimeIndex

- Can use time as an **index**
- data = [('2017 11 30', 48)]('2017 - 12 - 02', 45),(2017 - 12 - 03', 44),(2017 - 12 - 04', 48)dates, temps = zip(\*data)

s = pd.Series(temps, pd.to datetime(dates))

- Accessing a particular time or checking equivalence allows any string that can be interpreted as a date:
  - s['12/04/2017'] Or s['20171204']
- Using a less specific string will get all matching data:
  - s['2017-12'] returns the three December entries







### DatetimeIndex

• Time slices do not need to exist: - s['2017-12-01':'2017-12-31']









## Shifting Data

Leading or Lagging Data

```
In [95]: ts = Series(np.random.randn(4),
                    index=pd.date_range('1/1/2000', periods=4, freq='M'))
   • • • • •
In [96]: ts
                            In [97]: ts.shift(2)
                                                         In [98]: ts.shift(-2)
Out[96]:
                            Out[97]:
                                                         Out[98]:
2000-01-31
            -0.066748
                                                         2000-01-31
                            2000-01-31
                                               NaN
                                                                      -0.117388
            0.838639
                                               NaN
2000-02-29
                            2000-02-29
                                                         2000-02-29
                                                                      -0.517795
2000-03-31 -0.117388
                            2000-03-31 -0.066748
                                                         2000-03-31
                                                                            NaN
                                                                            NaN
2000-04-30 -0.517795
                            2000-04-30
                                          0.838639
                                                         2000-04-30
Freq: M, dtype: float64
                            Freq: M, dtype: float64
                                                         Freq: M, dtype: float64
```

• Shifting by time:

```
In [99]: ts.shift(2, freq='M')
Out[99]:
            -0.066748
2000-03-31
2000-04-30
           0.838639
            -0.117388
2000-05-31
2000-06-30 -0.517795
Freq: M, dtype: float64
```









### Shifting Time Series

• Data:

[(2017-11-30', 48), (2017-12-02', 45),('2017 - 12 - 03', 44), ('2017 - 12 - 04', 48)]

• Compute day-to-day difference in high temperature:

- s s.shift(1) (same as s.di
- 2017-11-30 NaN 2017 - 12 - 02 - 3.02017 - 12 - 03 - 1.04.0 2017 - 12 - 04









## Timedelta

- Compute differences between dates
- Lives in datetime module
- diff = parse date("1 Jan 2017") datetime.now().date() diff.days
- Also a pd. Timedelta object that take strings:
  - datetime.now().date() + pd.Timedelta("4 days")
- Also, Roll dates using anchored offsets from pandas.tseries.offsets import Day, MonthEnd

now = datetime(2011, 11, 17)In [107]: now + MonthEnd(2) Out[107]: Timestamp('2011-12-31 00:00:00')









### Time Zones

- Why?
- Coordinated Universal Time (UTC) is the standard time (basically equivalent to Greenwich Mean Time (GMT)
- Other time zones are UTC +/-a number in [1,12] • DeKalb is UTC-6 (aka US/Central); Daylight Saving Time is UTC-5









## Python, Pandas, and Time Zones

- Time series in pandas are time zone native
- The pytz module keeps track of all of the time zone parameters
  - even Daylight Savings Time
- Localize a timestamp using tz localize
  - -ts = pd.Timestamp("1 Dec 2016 12:30 PM")ts = ts.tz localize("US/Eastern")
- Convert a timestamp using tz\_convert
  - ts.tz convert ("Europe/Budapest")
- Operations involving timestamps from different time zones become UTC







### Frequency

- Generic time series in pandas are irregular
  - there is no fixed frequency
  - we don't necessarily have data for every day/hour/etc.
- Date ranges have frequency

```
In [76]: pd.date range(start='2012-04-01', periods=20)
Out[76]:
DatetimeIndex(['2012-04-01', '2012-04-02', '2012-04-03', '2012-04-04',
               '2012-04-05', '2012-04-06', '2012-04-07', '2012-04-08',
               '2012-04-09', '2012-04-10', '2012-04-11', '2012-04-12',
               '2012-04-13', '2012-04-14', '2012-04-15', '2012-04-16',
               '2012-04-17', '2012-04-18', '2012-04-19', '2012-04-20'],
              dtype='datetime64[ns]', freq='D')
```







## Lots of Frequencies (not comprehensive)

Alias	Offset type	Description
D	Day	Calendar daily
В	BusinessDay	Business daily
н	Hour	Hourly
T or min	Minute	Minutely
S	Second	Secondly
Lorms	Milli	Millisecond (1/1,000 of 1 second)
U	Місго	Microsecond (1/1,000,000 of 1 second)
Μ	MonthEnd	Last calendar day of month
BM	BusinessMonthEnd	Last business day (weekday) of month
MS	MonthBegin	First calendar day of month
BMS	BusinessMonthBegin	First weekday of month
W-MON, W-TUE,	Week	Weekly on given day of week (MON, TUE, WED, THU, FRI, SAT, or SUN)
WOM-1MON, WOM-2MON,	WeekOfMonth	Generate weekly dates in the first, second, third, or fourth week of the month (e.g., WOM-3FRI for the third Friday of each month)
Q-JAN, Q-FEB,	QuarterEnd	Quarterly dates anchored on last calendar day of each month, for year ending in indicated month (JAN, FEB, MAR, APR, MAY, JUN, JUL, AUG, SEP, OCT, NOV, or DEC
BQ-JAN, BQ-FEB,	BusinessQuarterEnd	Quarterly dates anchored on last weekday day of each month, for year ending in indicated month
QS-JAN, QS-FEB,	QuarterBegin	Quarterly dates anchored on first calendar day of each month, for year ending in indicated month
BQS-JAN, BQS-FEB,	BusinessQuarterBegin	Quarterly dates anchored on first weekday day of each month, for year ending in indicated month
A-JAN, A-FEB,	YearEnd	Annual dates anchored on last calendar day of given month (JAN, FEB, MAR, APR, MAY, JUN, JUL, AUG, SEP, OCT, NOV, or DEC)
BA-JAN, BA-FEB,	BusinessYearEnd	Annual dates anchored on last weekday of given month
AS-JAN, AS-FEB,	YearBegin	Annual dates anchored on first day of given month
BAS-JAN, BAS-FEB,	BusinessYearBegin	Annual dates anchored on first weekday of given month

### D. Koop, CSCI 640/490, Spring 2023

### [W. McKinney, Python for Data Analysis]



Northern Illinois University









## Resampling

- Could be
  - downsample: higher frequency to lower frequency
  - upsample: lower frequency to higher frequency
  - neither: e.g. Wednesdays to Fridays
- resample method: e.g. ts.resample('M').mean()

Argument	Description
freq	String or DateOffset indicating desired resampled
axis	Axis to resample on; default axis=0
fill_method	How to interpolate when upsampling, as in 'ff
closed	In downsampling, which end of each interval is c
label	In downsampling, how to label the aggregated re 9:30 to 9:35 five-minute interval could be labeled
loffset	Time adjustment to the bin labels, such as '-1s second earlier
limit	When forward or backward filling, the maximum
kind	Aggregate to periods ( 'period ' ) or timestamp time series has
convention	When resampling periods, the convention ('state to high frequency; defaults to 'end'

### D. Koop, CSCI 640/490, Spring 2023

frequency (e.g., 'M', '5min', or Second(15))

Fill' or 'bfill'; by default does no interpolation closed (inclusive), 'right' or 'left' result, with the 'right' or 'left' bin edge (e.g., the

d 9:30 or 9:35)

' / Second(-1) to shift the aggregate labels one

number of periods to fill

ps ('timestamp'); defaults to the type of index the

art' or 'end') for converting the low-frequency period

### [W. McKinney, Python for Data Analysis]









## Downsampling

- Need to define bin edges which are used to group the time series into intervals that can be aggregated
- Remember:
  - Which side of the interval is closed
  - How to label the aggregated bin (start or end of interval)



	9:02	9:03	9:04	9:05			
_							
	9:02	9:03	9:04	9:05			
				Ť			
	label=ˈrigh						







## Upsampling

### No aggregation necessary

In [222]: f	rame							
Out[222]:								
	Colorado	Texas	New York	Ohio				
2000-01-05	-0.896431	0.677263	0.036503	0.087102				
2000-01-12	-0.046662	0.927238	0.482284	-0.867130				
					In [225]: f	frame.resam	ple('D').f	fill()
In [223]:	df dailv =	frame.resa	<pre>ample('D')</pre>	.asfreg()	Out[225]:			
[ ] •	<u> </u>					Colorado	Texas	New Y
In [224]:	df dailv				2000-01-05	-0.896431	0.677263	0.036
Out[224]:					2000-01-06	-0.896431	0.677263	0.036
	Colorado	Texas	New York	Ohio	2000-01-07	-0.896431	0.677263	0.036
2000-01-05	-0.896431	0.677263	0.036503	0.087102	2000-01-08	-0.896431	0.677263	0.036
2000-01-06	NaN	NaN	NaN	NaN	2000-01-09	-0.896431	0.677263	0.036
2000-01-07	NaN	NaN	NaN	NaN	2000-01-10	-0.896431	0.677263	0.036
2000-01-08	NaN	NaN	NaN	NaN	2000-01-11	-0.896431	0.677263	0.036
2000-01-09	NaN	NaN	NaN	NaN	2000-01-12	-0.046662	0.927238	0.482
2000-01-10	NaN	NaN	NaN	NaN				
2000-01-11	NaN	NaN	NaN	NaN				
2000-01-12	-0.046662	0.927238	0.482284	-0.867130				

### D. Koop, CSCI 640/490, Spring 2023



Texas New York

0.677263 0.036503 0.087102

0.677263 0.036503 0.087102

0.677263 0.036503 0.087102

0.677263 0.036503 0.087102

0.677263 0.036503 0.087102

0.677263 0.036503 0.087102

0.677263 0.036503 0.087102

0.927238 0.482284 -0.867130

Ohio







9	13	4	11	3	8	
---	----	---	----	---	---	--







7.8

9	13	4	11	3	8	
---	----	---	----	---	---	--







7.8







7.8 7.0











7.8 7.0 8.3







## Window Functions

- then slide that window ahead. Repeat.
- rolling: smooth out data
- Specify the window size in rolling, then an aggregation method
- Result is set to the right edge of window (change with center=True)
- Example:
  - df.rolling('180D').mean()
  - df.rolling('90D').sum()

### D. Koop, CSCI 640/490, Spring 2023

• Idea: want to aggregate over a window of time, calculate the answer, and







### Interpolation

- algorithms
- Apply after resample

### D. Koop, CSCI 640/490, Spring 2023

### • Fill in the missing values with computed best estimates using various types of









### Sales Data by Month







### Resampled Sales Data (ffill)







## Resampled with Linear Interpolation (Default)



### D. Koop, CSCI 640/490, Spring 2023





42

## Resampled with Cubic Interpolation







## Piecewise Cubic Hermite Interpolating Polynomial







## 90-Day Rolling Window (Mean)



![](_page_53_Picture_3.jpeg)

![](_page_53_Picture_5.jpeg)

## 180-Day Rolling Window (Mean)

![](_page_54_Figure_1.jpeg)

![](_page_54_Picture_3.jpeg)

![](_page_54_Picture_5.jpeg)

### Time Series Databases

- Most time series data is heavy **inserts**, few updates
- Also analysis tends to be on ordered data with trends, prediction, etc.
- Can also consider stream processing
- Focus on time series allows databases to specialize
- Examples:
  - InfluxDB (noSQL)
  - TimescaleDB (SQL-based)

![](_page_55_Picture_11.jpeg)

![](_page_55_Picture_13.jpeg)

## Time Series Database Motivation

- Boeing 787 produces 500GB sensor data per flight
- Purposes
  - IoT
  - Monitoring large industrial installations
  - Data analytics
- Metrics (regular) and Events (irregular)
- Events can be obtained from metrics via binning

![](_page_56_Picture_13.jpeg)

![](_page_56_Picture_15.jpeg)

## What is a Time Series Database?

- A DBMS is called TSDB if it can
  - store a row of data that consists of timestamp, value, and optional tags - store multiple rows of time series data grouped together

  - can query for rows of data
  - can contain a timestamp or a time range in a query

	ul1	1 "SELECT * FROM ul1 WHERE time >= '2016-07-12T12:10:00Z"						Ζ'"
	time	generated	message_subtype	scaler	short_id	tenant	value	
	2016-07-12T11:51:45Z	"true"	"34"	"4"	"3"	"saarlouis"	465110000	
	2016-07-12T11:51:45Z	"true"	"34"	"-6"	"2"	"saarlouis"	0.06196699999999994	
$- \langle$	2016-07-12T12:10:00Z	"true"	"34"	"7"	"5"	"saarlouis"	4937000000	
	2016-07-12T12:10:00Z	"true"	"34"	"6"	"2"	"saarlouis"	18573000000	
	2016-07-12T12:10:00Z	"true"	"34"	"5"	"7"	"saarlouis"	5902300000	

### D. Koop, CSCI 640/490, Spring 2023

![](_page_57_Picture_9.jpeg)

![](_page_57_Picture_10.jpeg)

![](_page_57_Picture_12.jpeg)

49

## Storing Time Series Data in a RDBMS

- Timestamp as a primary key
- Tags and timestamp as combined primary key • Use an auto-incrementing primary key (timestamp is a normal attribute)

![](_page_58_Picture_5.jpeg)

![](_page_58_Picture_6.jpeg)

![](_page_58_Picture_8.jpeg)

![](_page_58_Picture_9.jpeg)

## Gorilla Motivation

- Large-scale internet services rely on lots of services and machines
- Want to monitor the health of the systems
- Writes dominate
- Want to detect state transitions
- Must be highly available and fault tolerant

![](_page_59_Figure_10.jpeg)

![](_page_59_Picture_11.jpeg)

![](_page_59_Picture_12.jpeg)

![](_page_59_Picture_14.jpeg)

![](_page_59_Picture_15.jpeg)

![](_page_59_Picture_16.jpeg)

## Gorilla Requirements

- 2 billion unique time series identified by a string key. 700 million data points (time stamp and value) added per minute.
- Store data for 26 hours.
- More than 40,000 queries per second at peak.
- Reads succeed in under one millisecond.
- Support time series with 15 second granularity (4 pts/minute per time series).
- Two in-memory, not co-located replicas (for disaster recovery capacity).
- Always serve reads even when a single server crashes.
- Ability to quickly scan over all in memory data.
- Support at least 2x growth per year.

![](_page_60_Picture_12.jpeg)

![](_page_60_Picture_13.jpeg)

![](_page_60_Picture_15.jpeg)

![](_page_60_Picture_16.jpeg)

![](_page_60_Picture_17.jpeg)

### Gorilla

- In-memory DB
- Data: 3-tuple string key, 64-bit timestamp integer, double-precision float
- Integer compression didn't work

![](_page_61_Picture_6.jpeg)

![](_page_61_Picture_8.jpeg)

![](_page_61_Picture_9.jpeg)

![](_page_61_Picture_10.jpeg)

### Time Series Data Patterns

![](_page_62_Figure_1.jpeg)

- Numerical Data Features:
  - Scale
  - Delta -
  - Repeat
  - Increase
- Text Data Features
  - Value
  - Character

![](_page_62_Picture_11.jpeg)

![](_page_62_Picture_12.jpeg)

![](_page_62_Picture_14.jpeg)

![](_page_62_Picture_15.jpeg)

## Gorilla Compression

![](_page_63_Figure_1.jpeg)

### D. Koop, CSCI 640/490, Spring 2023

### 

Northern Illinois University

![](_page_63_Picture_5.jpeg)

![](_page_63_Picture_6.jpeg)

![](_page_63_Picture_7.jpeg)

### Delta of Delta Compression

- Data usually recorded at regular intervals
- Deltas: 60, 60, 59, 61
- Delta of deltas (D): 0, -1, 2
- Variable-length encoding:
  - $D = 0 \rightarrow 0$
  - D in [-63,64]  $\rightarrow$  10 + value (7 bits)
  - D in  $[-255, 256] \rightarrow 110 + value (9 bits)$
  - D in  $[-2047, 2048] \rightarrow 1110 + value (12 bits)$
  - else  $\rightarrow$  1111 + value (32 bits)
- 1 bit 96% of the time

![](_page_64_Figure_13.jpeg)

[Pelkonen et al., 2015]

![](_page_64_Picture_15.jpeg)

![](_page_64_Picture_17.jpeg)

![](_page_64_Picture_18.jpeg)

![](_page_64_Picture_19.jpeg)

![](_page_64_Picture_20.jpeg)

## XOR Representation

Decimal	Double Representation	XOR with previous
12	0x40280000000000000	
24	0x40380000000000000	0x0010000000000000
15	0x402e0000000000000	0x0016000000000000
12	0x40280000000000000	0x000000000000000000000000000000000000
35	0x4041800000000000	0x00 <mark>6980</mark> 0000000000
Decimal	<b>Double Representation</b>	XOR with previous
15.5	0x402f0000000000000	
14.0625	0x402c200000000000	0x0003200000000000
3.25	0x400a000000000000	0x0 <mark>0</mark> 26200000000000
8.625	0x4021400000000000	0x002b400000000000
13.1	0x402a33333333333333	0x000b73333333333333

- Values usually do not change significantly Look at XOR
  - Same  $\rightarrow 0$
  - Changes in Meaningful Bits
    - Same as previous value  $\rightarrow$  10 + changed bits
    - Outside previous value  $\rightarrow$  11 + leading zeros + length of meaningful bits + bits

![](_page_65_Picture_9.jpeg)

![](_page_65_Picture_11.jpeg)

![](_page_65_Picture_12.jpeg)

![](_page_65_Picture_13.jpeg)

## Enabling Gorilla Features

- Correlation Engine: "What happened around the time my service broke?"
- Charting: Horizon charts to see outliers and anomalies
- Aggregations: Rollups locally in Gorilla every couple of hours

![](_page_66_Figure_6.jpeg)

![](_page_66_Picture_7.jpeg)

![](_page_66_Picture_8.jpeg)

![](_page_66_Picture_10.jpeg)

![](_page_66_Picture_11.jpeg)

![](_page_66_Picture_12.jpeg)

![](_page_66_Picture_13.jpeg)

## Gorilla Lessons Learned

- Prioritize recent data over historical data
- Read latency matters
- High availability trumps resource efficiency
  - Withstand single-node failures and "disaster events" that affect region
  - "[B]uilding a reliable, fault tolerant system was the most time consuming part of the project"
  - "[K]eep two redundant copies of data in memory"

![](_page_67_Picture_9.jpeg)

![](_page_67_Picture_10.jpeg)

![](_page_67_Picture_12.jpeg)

![](_page_67_Picture_13.jpeg)

![](_page_67_Picture_14.jpeg)