

# consortium

Part of the new  
R/Insurance  
Webinar Series

High performance  
programming in R  
(arrow package and Parquet files)

31 January 2024

# Welcome to the webinar!

R/insurance webinar series

- 1) From Excel to programming in R
- 2) From programming in R to putting R into production
- 3) R performance culture
- 4) High performance programming in R**

Delivered on behalf of the R Consortium by Georgios Bakoloukas and Benedikt Schamberger, Actuarial Control, Group Risk Management, Swiss Re

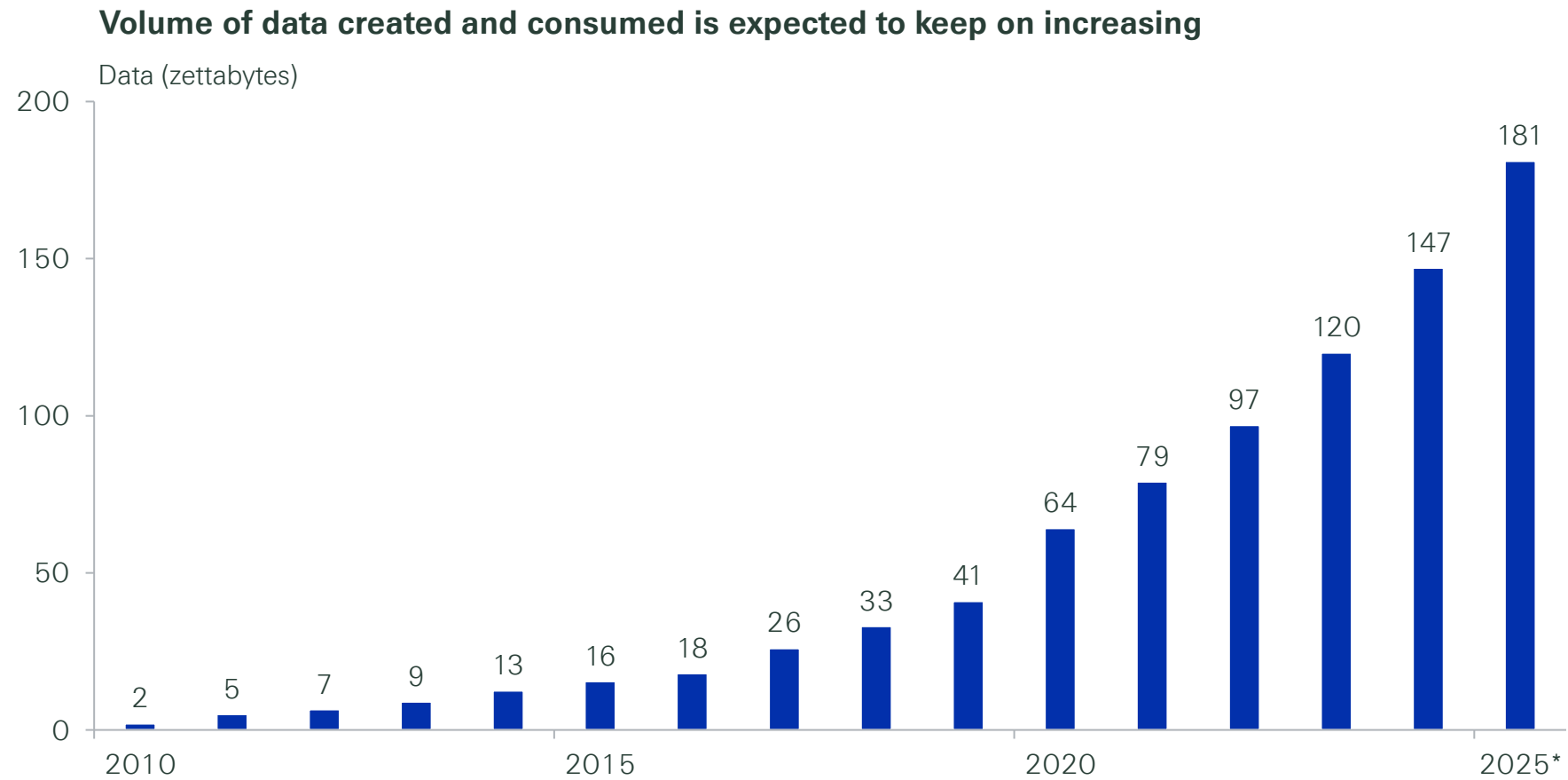
## Background to Swiss Re's R community

Large actuarial R programming, Atelier, community

- Swiss Re internal R community sponsored by our Group Chief Actuary Philip Long ([Atelier programme](#))
- 2000+ community with 500+ regular coders who also support each other
- The case we see today relates to code optimisations we did for experience study work, ie comparing how an insurance portfolio performed to initial expectations
- Views expressed belong solely to the speakers and not necessarily to the speaker's employer

# Data sizes increasing

Sizes have been growing exponentially and are expected to keep this trend



Source: [Statista](#): Volume of data/information created, captured, copied, and consumed worldwide from 2010 to 2020, with forecasts from 2021 to 2025, September 2022

# Parquet file format

Uses common compression methods like RLE, dictionaries and other optimisations



Fast and efficient columnar-storage

Run length encoding (RLE)		Dictionary encoding		Performance optimisation																																
<table border="1"> <thead> <tr> <th>Row</th> <th>Name,...</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>Benedikt,...</td> </tr> <tr> <td>2</td> <td>Benedikt,...</td> </tr> <tr> <td>...</td> <td>...</td> </tr> <tr> <td>1,000,000</td> <td>Benedikt,...</td> </tr> </tbody> </table>	Row	Name,...	1	Benedikt,...	2	Benedikt,...	...	...	1,000,000	Benedikt,...	<table border="1"> <thead> <tr> <th>Item</th> <th>Times</th> </tr> </thead> <tbody> <tr> <td>Benedikt</td> <td>1,000,000</td> </tr> </tbody> </table>	Item	Times	Benedikt	1,000,000	<p><b>Company name</b> Swiss Reinsurance Company, Ltd.</p>	<p><b>Company name</b> 42</p>	<table border="1"> <thead> <tr> <th>Row group</th> <th>Age</th> <th>Age &lt; 20</th> </tr> </thead> <tbody> <tr> <td rowspan="3">1</td> <td>30</td> <td rowspan="3">min: 27 max: 31 <b>min &gt; 20 Skip</b></td> </tr> <tr> <td>31</td> </tr> <tr> <td>27</td> </tr> <tr> <td rowspan="3">2</td> <td>19</td> <td rowspan="3">min: 14 max: 19 <b>max &lt; 20 Search</b></td> </tr> <tr> <td>14</td> </tr> <tr> <td>17</td> </tr> <tr> <td rowspan="3">3</td> <td>55</td> <td rowspan="3">min: 23 max: 55 <b>min &gt; 20 Skip</b></td> </tr> <tr> <td>41</td> </tr> <tr> <td>23</td> </tr> </tbody> </table>	Row group	Age	Age < 20	1	30	min: 27 max: 31 <b>min &gt; 20 Skip</b>	31	27	2	19	min: 14 max: 19 <b>max &lt; 20 Search</b>	14	17	3	55	min: 23 max: 55 <b>min &gt; 20 Skip</b>	41	23
Row	Name,...																																			
1	Benedikt,...																																			
2	Benedikt,...																																			
...	...																																			
1,000,000	Benedikt,...																																			
Item	Times																																			
Benedikt	1,000,000																																			
Row group	Age	Age < 20																																		
1	30	min: 27 max: 31 <b>min &gt; 20 Skip</b>																																		
	31																																			
	27																																			
2	19	min: 14 max: 19 <b>max &lt; 20 Search</b>																																		
	14																																			
	17																																			
3	55	min: 23 max: 55 <b>min &gt; 20 Skip</b>																																		
	41																																			
	23																																			
<p><b>Size</b> ~8,000,000 B</p>	<p><b>Size</b> ~16 B</p>	<p><b>Size</b> ~31 B</p>	<p><b>Size</b> ~4 B</p>																																	

Examples for illustration, exact sizes and optimisations will depend on your data and options used in creating the parquet file

# Moving from CSV to Parquet with arrow

Drop-in replacement for CSV

CSV (with data.table)



Parquet (with arrow)



## Reading data

```
my_data <- data.table::fread("my_data.csv")
```



```
my_data <- arrow::read_parquet("my_data.parquet")
```

## Writing data

```
data.table::fwrite(my_data, file = "my_data.csv")
```



```
arrow::write_parquet(my_data, sink = "my_data.parquet")
```

# Case study: Aggregating exposure information

Regular dplyr and partitioned parquet

## Regular dplyr



```
df <- fread("my_data.csv") # or read_parquet("my_data.parquet")

df |>
  filter(calendar_year == 2019) |>
  summarise(
    total_exposure = sum(exposure_amount),
    .by = c("benefit_id", "insured_age_group", "region")
  )
```



## Partitioned parquet



```
df <- open_dataset("my_data") # partitioned parquet folder

df |>
  filter(calendar_year == 2019) |>
  summarise(
    total_exposure = sum(exposure_amount),
    .by = c("benefit_id", "insured_age_group", "region")
  ) |>
  collect()
```

# Why switching to Parquet is worthwhile

Smaller files, faster reading and potentially much faster aggregations

	File size	Read time		Grouping	
		data.frame	arrow Table	Single file Table	Partitioned files
<b>CSV</b>	<b>9.3GB</b>	<b>50s</b>	<b>50s</b>	<b>3.5s</b>	<b>3.5s</b>
	$\sim \frac{1}{5}$ th	$\sim \frac{1}{10}$ th	$\sim \frac{1}{25}$ th	$\sim \frac{1}{6}$ th	$\sim \frac{1}{15}$ th
<b>Parquet</b>	<b>1.9GB</b>	<b>5s</b>	<b>2s</b>	<b>0.6s</b>	<b>0.2s</b>

Test file: L&H insurance treaty about 24m rows, with 51 columns, size 9.3GB (CSV)



## Parquet with arrow is a valuable tool

Give it a try, you can switch back at any time

Parquet has  
been hugely  
successful  
worldwide

It is typically  
faster, and can  
be much faster  
if used right

Easy to use in  
R, but may  
present  
challenges in  
Excel

## R Consortium Impact

- R Consortium Community **Grants** and Sponsorships Over **USD \$1.4 Million**
- Organize large scale **collaborative projects**
  - R Validation Hub
  - R-Ladies
  - Diversity and Inclusion Working Group
- Co-host multidisciplinary **data science forums**
  - Stanford Data Institute
- Direct support for key **R events**
  - R/Medicine, R/Pharma, useR!, LatinR, more
- Direct support for **R User Groups**



**Organizations Can  
Become a Member  
Today!**

Email Joseph Rickert at  
**director@r-  
consortium.org**  
to set up first call