

DELPHIX

EE. hroug

17-20.10.2017

# Efficient Test Data Management

**Marcin Przepiorowski** | Senior Technical Principal | October, 2017

# About me

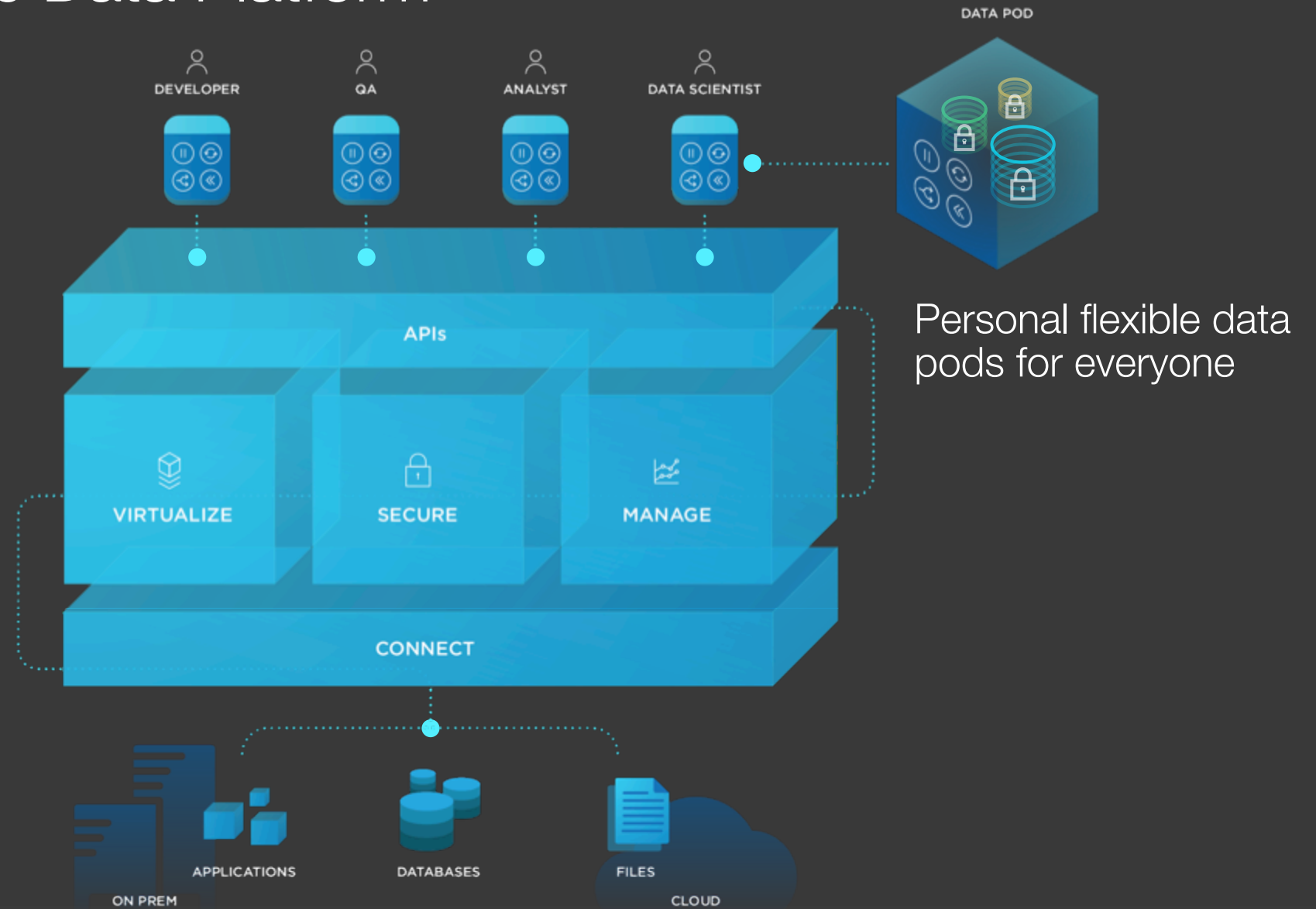
Oracle consultant/DBA since 2000

co-developer of OraSASH – free ASH/AWR like repository

Blogger ???



# Delphix Dynamic Data Platform



This session is focused on the tools and processes.

No actual database or vendor platform special knowledge is required to gain value from the session.



# Real life examples

Test Data Management - Tools and processes  
Security

# Example 1



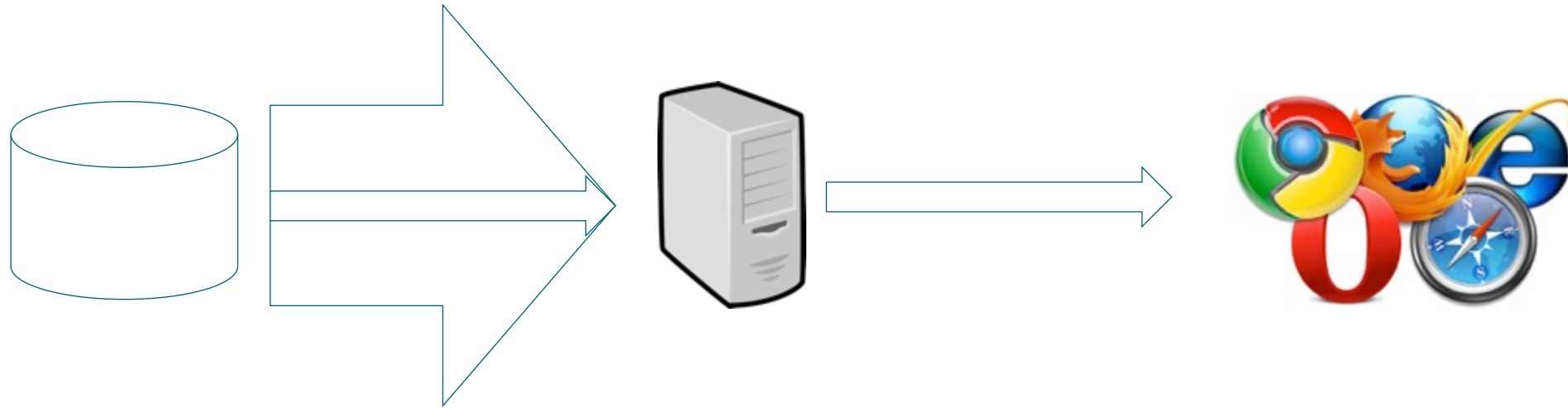
Test data are not like a whiskey or wine  
keeping them long doesn't increase a value

- Migration project based on development environment created 2 years before project kick off
- Risk to production migration due to data changes

[https://cdn.pixabay.com/photo/2016/03/31/15/23/cabinet\\_1293245\\_1280.png](https://cdn.pixabay.com/photo/2016/03/31/15/23/cabinet_1293245_1280.png) <https://www.flickr.com/photos/dionhinchcliffe/9505044956>



## Example 2



- Development project done for 18 months without tests on even 1 % of the production data
- Major issues found to close to the scheduled go – live

# Test Data Management

Test Data management is very critical during the test life cycle.

Over 80% of organizations stated that RECEIVING or REFRESHING the data to perform tests was the largest consumer of testing time, (over 90%) leaving actual work to consume less than 10% of the overall testing scenario.

[https://www.tutorialspoint.com/software\\_testing\\_dictionary/test\\_data\\_management.htm](https://www.tutorialspoint.com/software_testing_dictionary/test_data_management.htm)



Real life examples



# Test Data Management - Tools

Security

# What is a good TDM process

A process that assists in delivering a data sets for testing / development on time.

- Delivers a “right size” and secure datasets on time
- Has an ability to quickly isolate and deliver test cases for development to investigate.
- Has an ability to identify code and data changes by versions

# Test Data Management



# Test data

| METHODS                                | PROS   | CONS  |
|--|--|---|
| <b>Cloning</b>                         | Relatively simple to   | <ul style="list-style-type: none"> <li>Expensive in terms of hardware, license and support costs</li> <li>Time-consuming: Increases the time required to run test cases due to large data volumes</li> <li>Not agile: Developers, testers and QA staff can't refresh the test data</li> <li>Inefficient: Developers and testers can't create targeted test data sets for specific test cases or validate data after runs</li> <li>Requires collaborative between DBA and testing teams</li> <li>Not scalable across multiple data sources or applications</li> <li>Laborious: Production systems are typically large</li> <li>Risky: Access to real company data is required (developers, testers and QA staff need valid business reason to access sensitive data such as customer information)</li> </ul> |
| <b>Generating synthetic test data</b>  | Safe   | <ul style="list-style-type: none"> <li>Requires highly skilled DBAs with deep knowledge of the underlying data relationships that might not be formally detailed in the schema</li> <li>Tedious: DBAs often manually include errors and set boundary conditions within the synthetic data set to ensure a robust testing process, which adds time to the test data creation process</li> <li>Challenging: Despite the time and effort put forth by the DBAs, it is challenging to work with because synthetic test data does not always retain the proper context</li> <li>Time-consuming: Process is slower and can be error-prone</li> </ul>  |
| <b>Subsetting production databases</b> | Less expensive compared to cloning or generating synthetic test data | <ul style="list-style-type: none"> <li>Skill-intensive: Without an automated solution, requires highly skilled DBAs and requires access to production data to protect sensitive data</li> </ul>   |

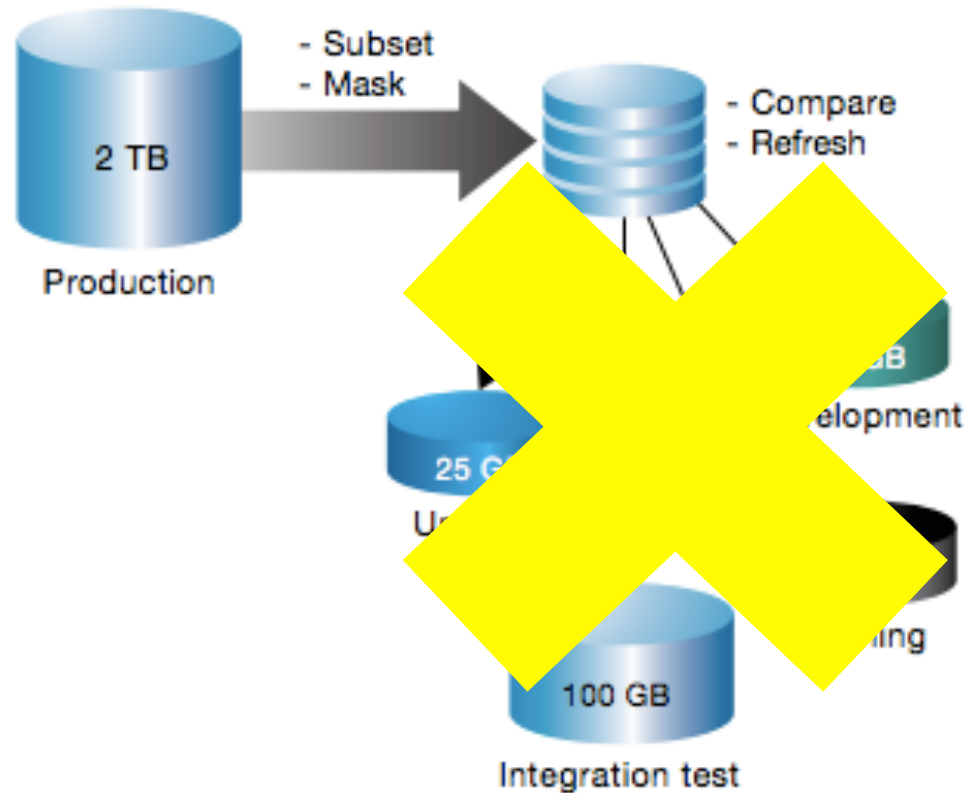
Time consuming

Resource consuming

Skills intensive

[http://www.informationweek.com/pdf\\_whitepapers/approved/1345732672\\_back\\_to\\_basics.pdf](http://www.informationweek.com/pdf_whitepapers/approved/1345732672_back_to_basics.pdf)

# What is the “Right Size”?



Development and test on “unrealistic” amounts of data can create a code quality issues.

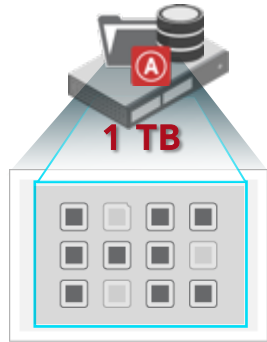
Production implementation can end up with scaling issues which are very costly to fix

# Sub-setting vs clone

- Complicated structures (ex. history of all objects in one table ) – not easy to subset
- Sub-setting requires a business logic to be implemented.
- Easiest option for sub-setting is if application has a proper archiving option

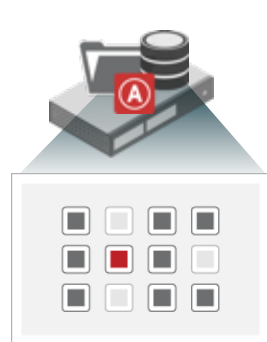
# Cloning issue

PRODUCTION  
Database/App Tier



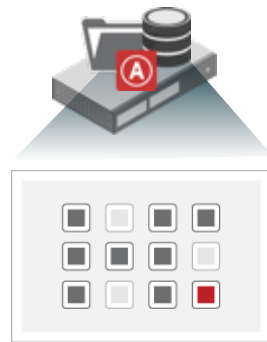
1 TB

QA



1 TB

TEST



1 TB

DEV



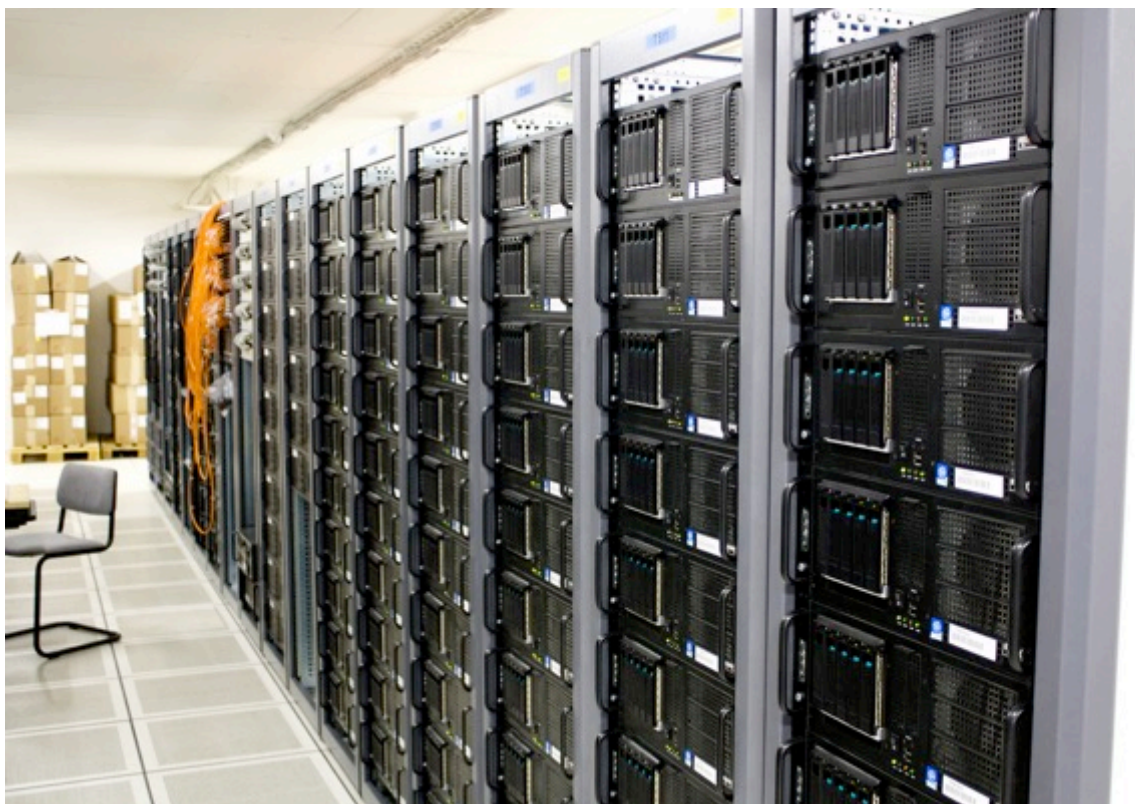
1 TB

PATCH TEST



1 TB





<https://www.flickr.com/photos/torkildr/3462607995>

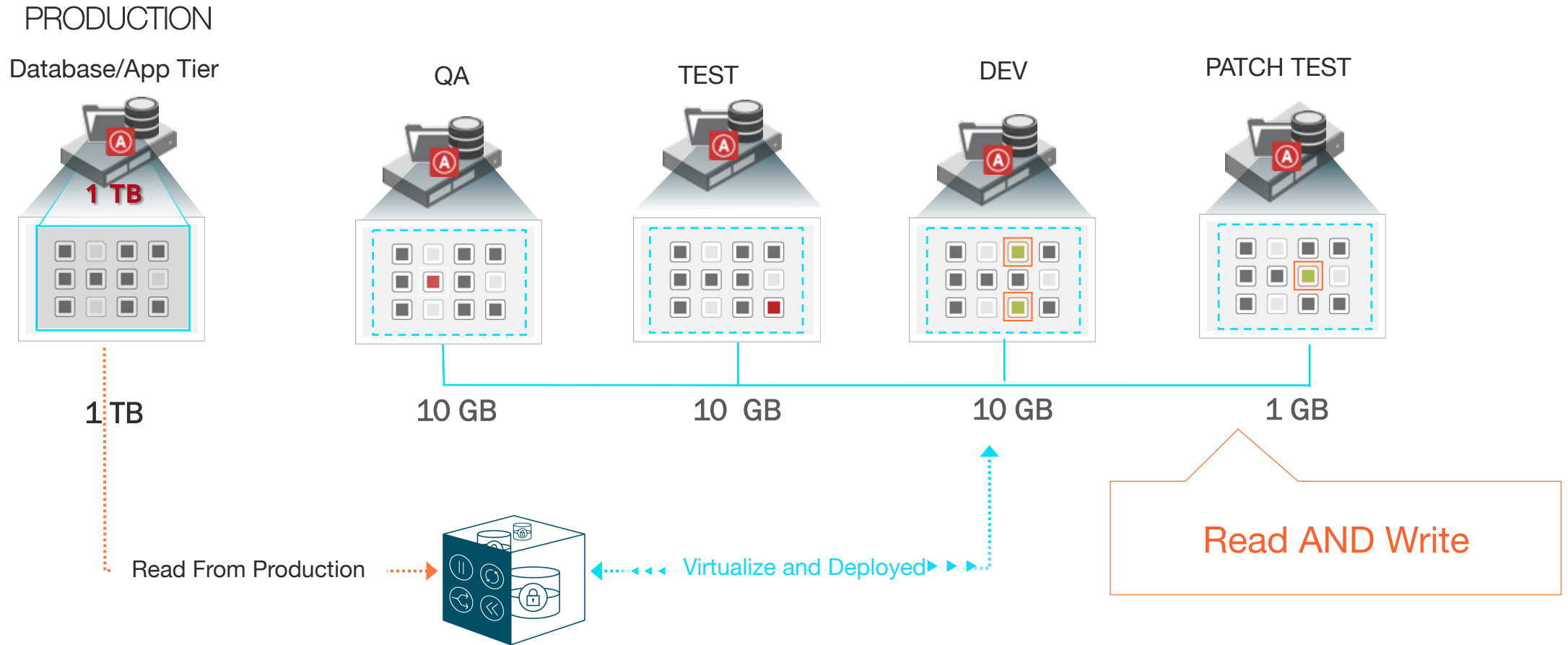


**ORACLE®**

VM



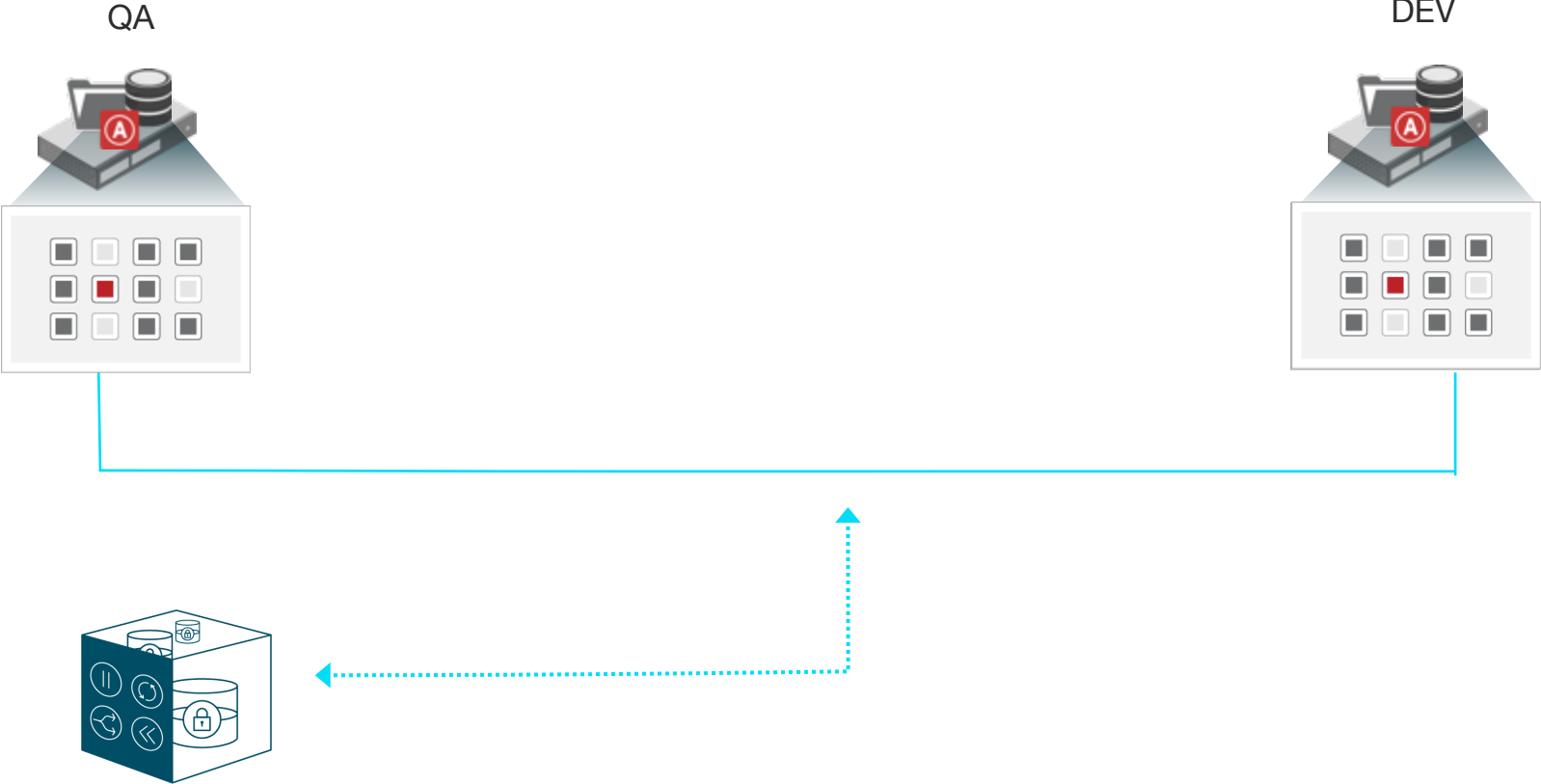
# Cloning issue - solution



# Isolate and deliver test data

- In case of data related issues, developers need to work on QA system to nail down the problem
- Potential tools and access issues

# Virtualization – sharing data between QA and Dev

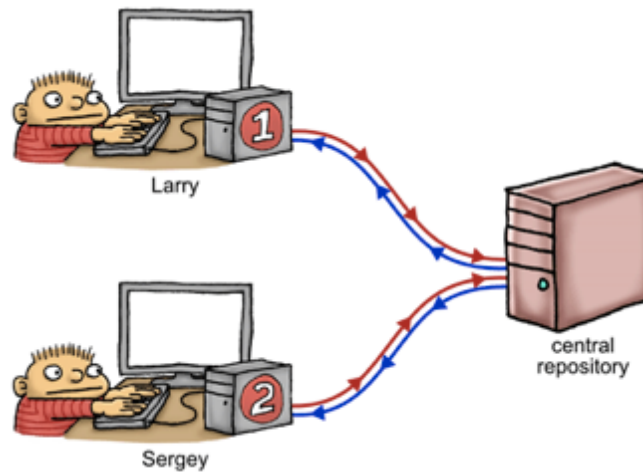


# Code and data changes / versions

- Most of RDBMSs are not good in keeping code and data versioning
- Branching a database is a challenge
- Rolling back changes is a challenge

# Code - Source Control

“A component of software configuration management, version control, also known as revision control or source control, is the management of changes to documents, computer programs, large web sites, and other collections of information.”



# Data – “Source” Control

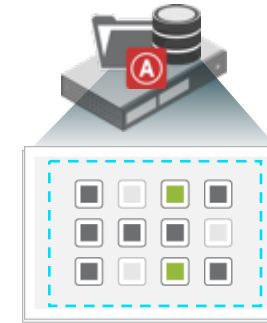
Ver. 1.7-Prod



28 Jan 2017



DEV

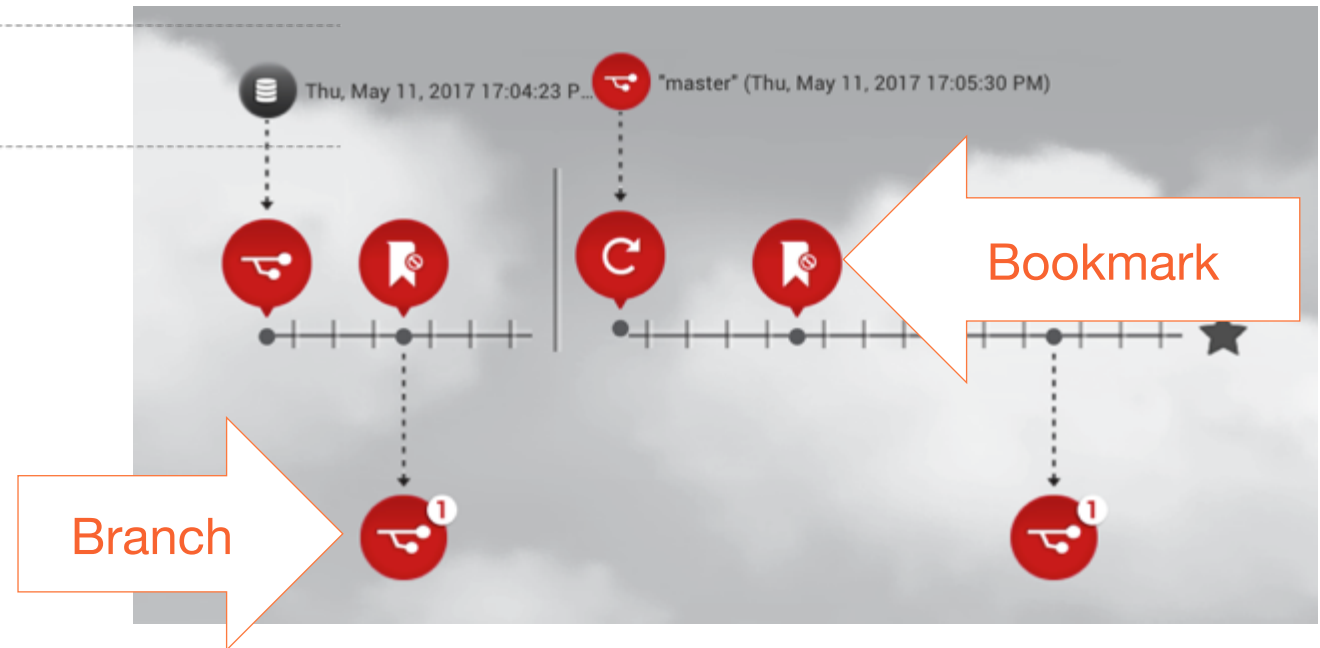


Ver. 2.0





# Data – “Source” Control



# Virtualization – How

- Storage based solutions – buy or DIY
- File system based solutions – buy or DIY
- Appliance solutions – buy

Real life examples



# Test Data Management - processes

Security



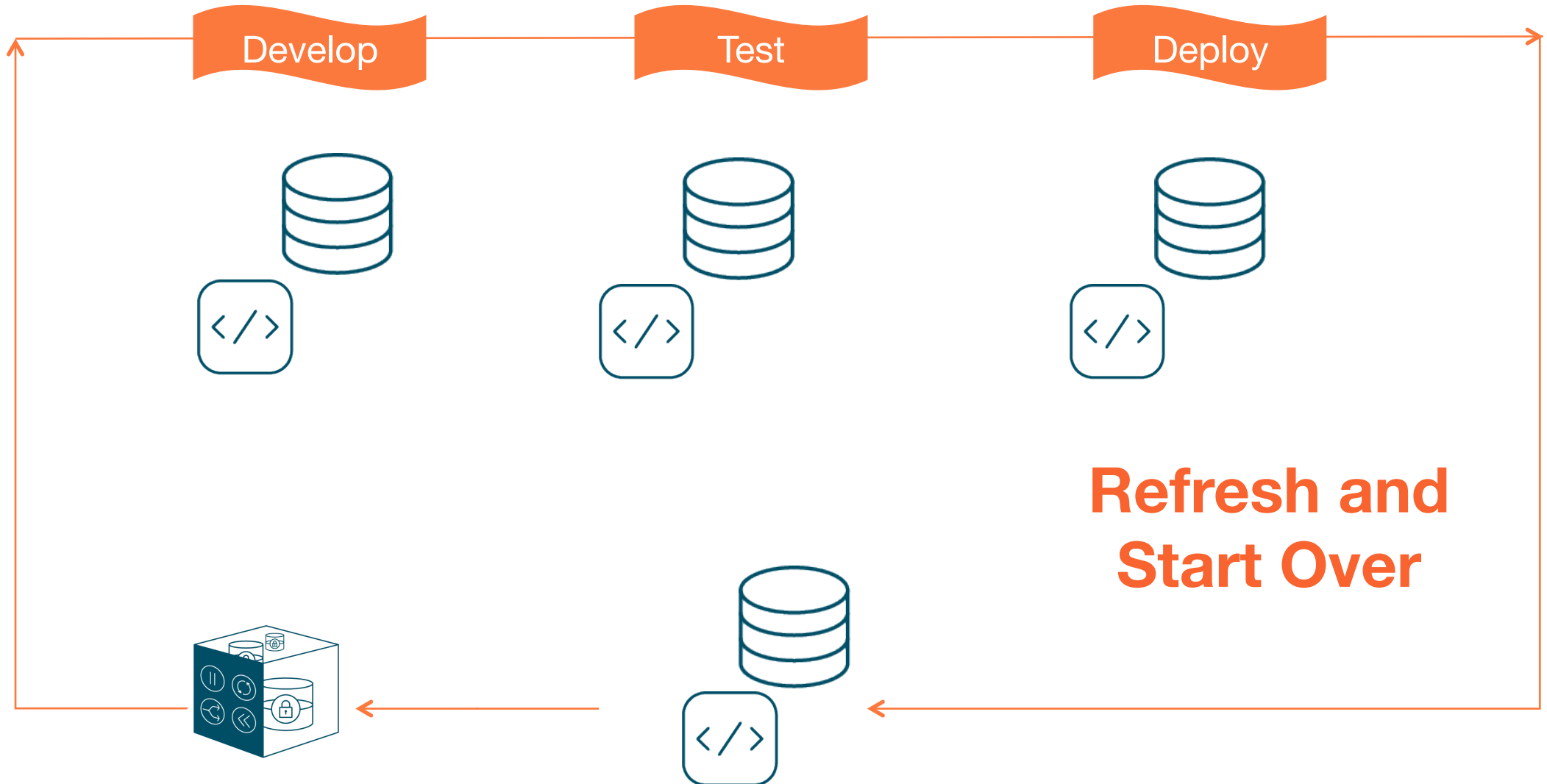
# Continuous ...



Data  
Virtualization



# Agile TDM



# Database Cloud Self Service Portal

## Create Database

Submit Cancel

### General

Service Template Thin Copy of CRM Production

\* Request Name SSA\_USER1 - Mon Sep 22 2014 16:14:52 EDT

\* Zone DbaaS\_Zone

\* Database SID CRMTest

\* Database Service Name CRMTest\_Srvc1

### Schedule Request

If Start Date is set to "Immediately", the timezone "Eastern Daylight Time (GMT -4:00)" will be used for End Date.

Start  Immediately  Later (UTC-05:00) New York - Eastern Time (ET)

Duration  Indefinitely  Until 9/22/2014 8:17:43 PM

### Deployment Input

\* User Name dbaas\_crm

\* User Password .....

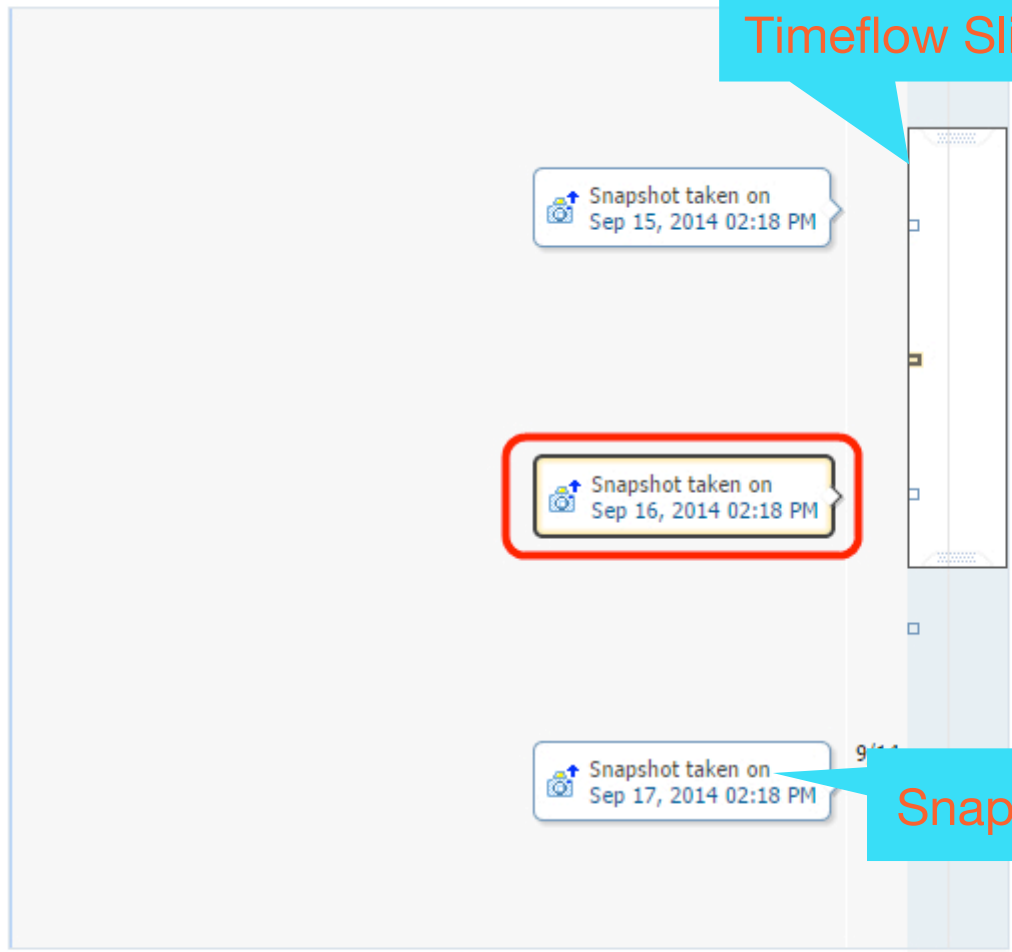
\* Confirm User Password .....

TIP The same password will be used for all schemas created as part of this request.

### Service Instance Properties

| Name     | Value |
|----------|-------|
| Optional |       |

## Snapshots



Timeflow Slider

Snapshot



### Datasets +

Filter: None Q



- Analytics
- targetcon  
CDB - Cannot contact source
- autofs  
vFiles - Stopped
- autotest**  
VDB - Running
- m... lsn  
base
- VDB - Stopped
- si4rac  
VDB - Running
- siclone  
VDB - Stopped
- vPDBtest  
VDB - Stopped
- Sources

Environments

### autotest ✎

Status **Timeflow** Configuration

**May 12, 2017 9:04 PM**  
 Europe/Dublin, GMT+0100  
**End Stamp** May 12, 2017 9:04 PM  
**Source Database** autotest  
**Database Version** 15.7 SP101  
**OS** Linux

**May 12, 2017 11:55 AM**  
 Europe/Dublin, GMT+0100  
**End Stamp** May 12, 2017 11:55 AM  
**Source Database** autotest  
**Database Version** 15.7 SP101  
**OS** Linux

Snapshot

Timeflow Slider



Selected Time Point May 12 2017 15 9 : 04 PM 🕒

Refresh VDB ↺ Rewind VDB ➡ V2P + ⌂

Actions

oracle sybasecont default

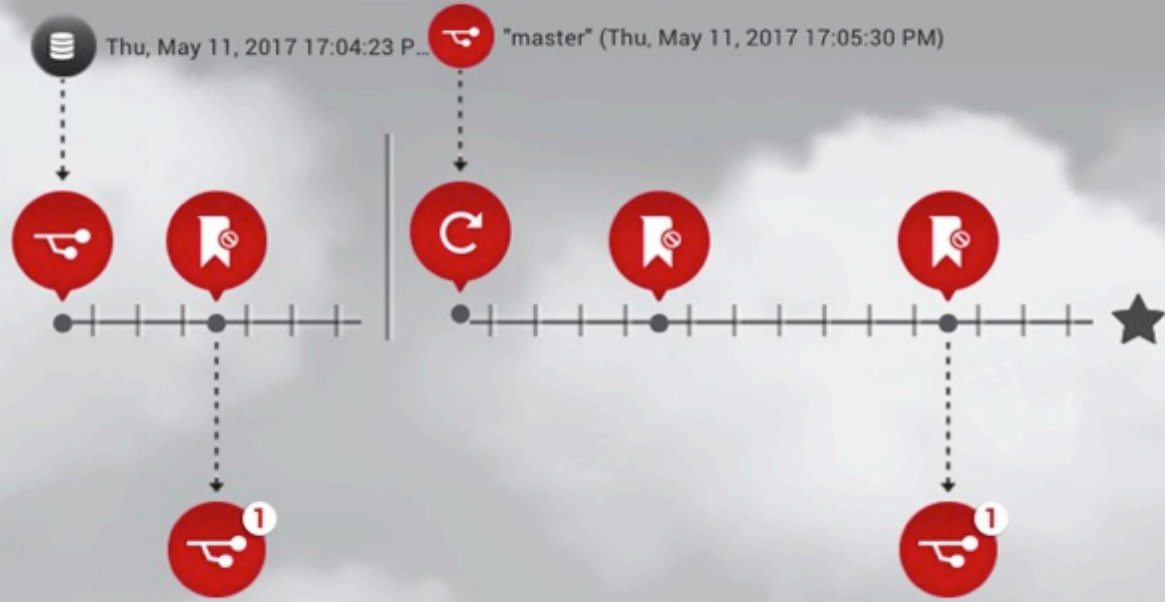
TIME

CONTAINERS

BRANCHES

BOOKMARKS

- default
- Version 2.0
- Branch 3.0
- New Branch 1



ACTIVATE

BOOKMARK

BRANCH

SHARE

REFRESH

RESTORE

RESET

STOP

LOCK

Real life examples

Test Data Management - processes

Security

# Data Security

|  |                             |                                 |
|--|-----------------------------|---------------------------------|
| <b>Totals for Category: Banking/Credit/Financial</b> | <b># of Breaches: 36</b>    | <b># of Records: 26,262</b>     |
|  | <b>% of Breaches: 4.3%</b>  | <b>%of Records: 0.1%</b>        |
| <hr/>  |                             |                                 |
| <b>Totals for Category: Business</b>                 | <b># of Breaches: 375</b>   | <b># of Records: 2,548,225</b>  |
|  | <b>% of Breaches: 44.4</b>  | <b>%of Records: 8.6%</b>        |
| <hr/>  |                             |                                 |
| <b>Totals for Category: Educational</b>              | <b># of Breaches: 74</b>    | <b># of Records: 489,376</b>    |
|  | <b>% of Breaches: 8.8%</b>  | <b>%of Records: 1.6%</b>        |
| <hr/>  |                             |                                 |
| <b>Totals for Category: Government/Military</b>      | <b># of Breaches: 58</b>    | <b># of Records: 12,300,322</b> |
|  | <b>% of Breaches: 6.9%</b>  | <b>%of Records: 41.3%</b>       |
| <hr/>  |                             |                                 |
| <b>Totals for Category: Medical/Healthcare</b>       | <b># of Breaches: 302</b>   | <b># of Records: 14,400,946</b> |
|  | <b>% of Breaches: 35.7</b>  | <b>%of Records: 48.4%</b>       |
| <hr/>  |                             |                                 |
| <b>Totals for All Categories:</b>                    | <b># of Breaches: 845</b>   | <b># of Records: 29,765,131</b> |
|  | <b>% of Breaches: 100.0</b> | <b>%of Records: 100.0%</b>      |

[http://www.idtheftcenter.org/images/breach/ITRCBreachReport\\_2016.pdf](http://www.idtheftcenter.org/images/breach/ITRCBreachReport_2016.pdf)

# Confidential data



# Do I Have to Mask Data?

| Type of Data           | Year Passed        | Ruling   |
|------------------------|--------------------|--|
| Data Masking in the EU | 2014               | <a href="#"><u>ARTICLE 29 DATA PROTECTION</u></a>                          |
| GDPR                   | 2016               | <a href="#"><u>Regulation (EU) 2016/679</u></a>                            |
| HIPAA                  | 1996               | <a href="#"><u>Health Insurance Portability and Accountability Act</u></a> |
| PCI                    | 2016,<br>(Updated) | <a href="#"><u>Payment Card Industry Standards</u></a>                     |
| PII                    |                    | <a href="#"><u>Personably Identifiable Information</u></a>                 |
| SOX                    | 2002               | <a href="#"><u>Sarbanes-Oxley Act</u></a>                                  |

# Masking in the Picture

As 80% of data in a company are copies, then 80% of data won't be subject to security like a production environment. Securing this data is not just a priority, but in many cases, subject to legal ramifications, (i.e. PCI/PII)



# Masking in the Security officer picture

## Masking Requirements

- Masking shouldn't be reversible
- Easy to audit
- Masking should be a simple, automated, repeatable process

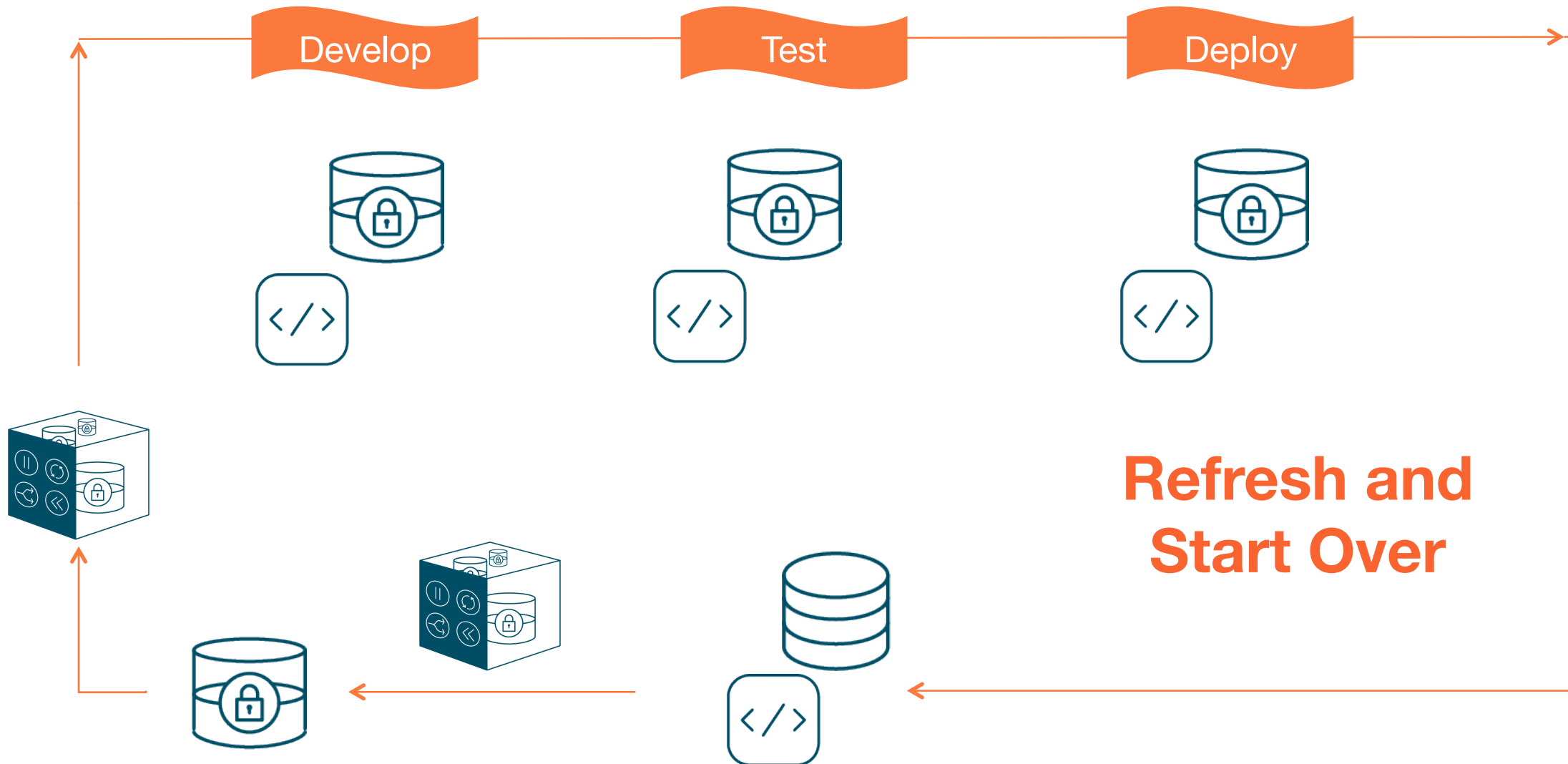
# Masking in the DBA picture

- Your data has to look same for optimizer after masking
  - Ex. age of customers
  - Ex. date of purchase
  - Ex. addresses
- Keep in mind correlation between data
- Keep in mind referential integrity

# How This All Comes Together...

- Virtualization is the key to fast, efficient and FULL copies of production environments for agile and automated testing for agile shops.
- Data masking that can be done once, easily maintained with a repeatable process via a strong discovery and implementation as part of the virtualization process secures the 80% of data that is outside the control of production.
- Virtualized environments that are built with development and testing in an Agile or DevOps environments makes it simple to accomplish what may see impossible and do so at light speed.

# Agile TDM



Marcin Przepiorowski  
Senior Technical Principal  
marcin@delphix.com  
@pioro