

Distributed Database Architecture for GDPR

Karthik Ranganathan
Co-Founder & CTO
Mar, 2019

Introduction



Karthik Ranganathan

Co-Founder & CTO, YugaByte
Nutanix Facebook Microsoft
IIT-Madras, University of Texas-Austin



@karthikr

What is YugaByte DB?



Distributed SQL (PostgreSQL compatible)

Highly scalable, Fault-Tolerance



High Performance

Low Latency Queries



Cloud Native

Multi-Cloud, multi-region and Kubernetes Ready

WHAT IS GDPR?

GDPR : General Data Protection Regulation

Persons in EU can control sharing and protection
of their personal data by businesses.



Personal Data, similar to

PII (Personally Identifiable Information)

- User name
- Email address
- Date of birth
- Bank details
- Location details
- Computer IP address

Control over personal data

- Consent & data location
- Data privacy and safety
- Right to be forgotten
- Data access on demand

Database concerns

- Notify on data breach
- Data portability
- Ability to fix errors in data
- Restrict processing

Application concerns

#1 USER CONSENT AND DATA LOCATION

Data must be stored in EU by default. Businesses need explicit user consent to move it outside.



Why is this hard?

- EU user data lives in that region
- Other countries have compliance regulation – more geo's
- Public clouds may not have coverage – hybrid deployments
- Architecture depends on data – multiple per service

Think Global Deployments first!

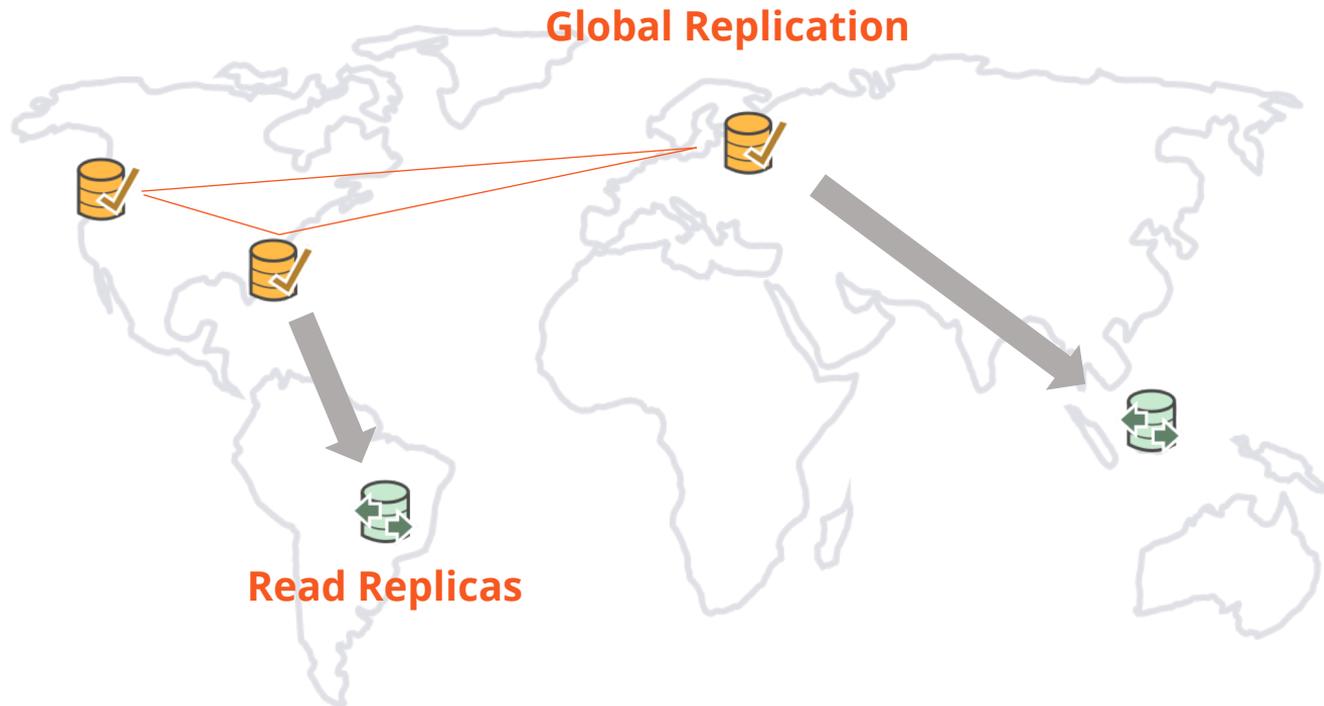
Example – online ecommerce site

- **Products** table needs globally replication – not PII data



Non-PII Data

Global Replication
with YugaByte DB

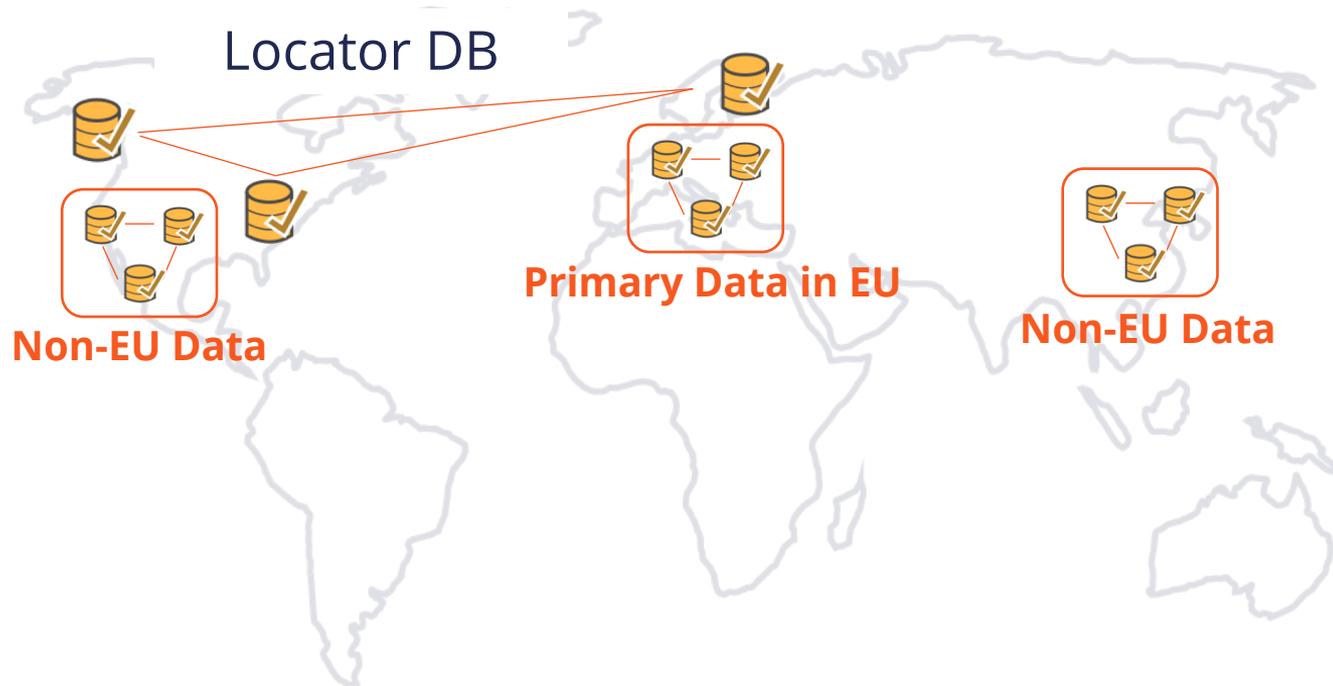


Example – online ecommerce site

- **Users, orders** and **shipments** needs locality – PII data
- **Product locations** table needs scale – may be PII

PII Data

Geo-Partitioning with YugaByte DB



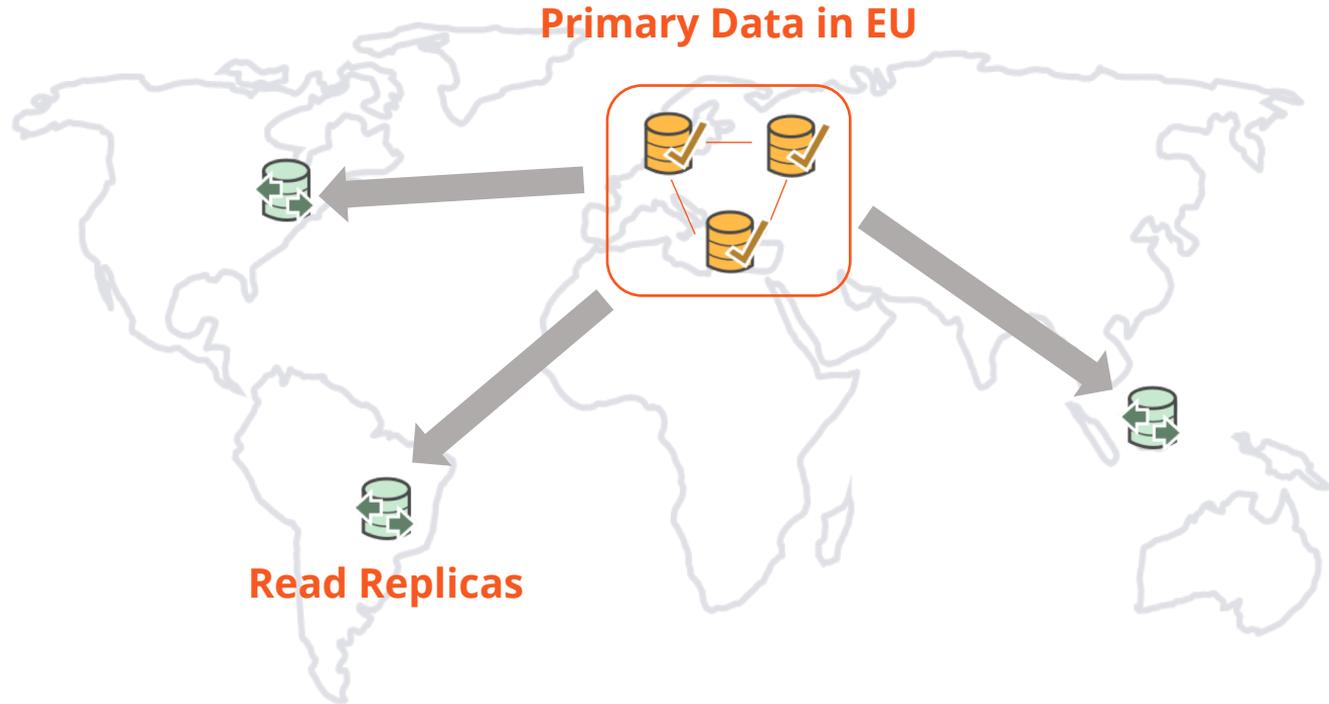
Replicate data on demand to other geo's

- User may be ok with replicating data
- Read replicas on demand (for remote, low-latency reads)
- Change data capture (for analytics)



PII Data with YugaByte DB

Read Replicas with
YugaByte DB



#2 DATA PRIVACY AND SAFETY

Data must be secured by using best practices by default. Users need to be notified on breach.

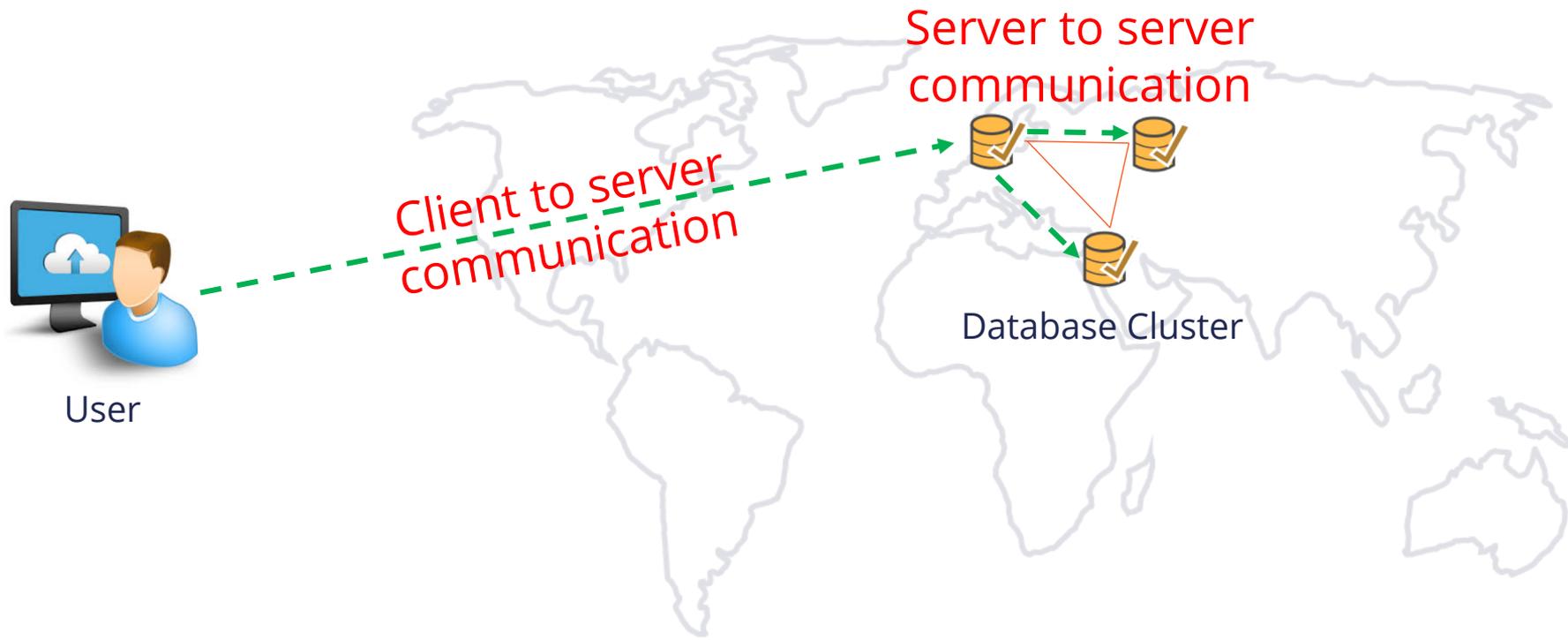


Implement end-to-end encryption on day #1

Encrypt All Network Communication

- Use TLS Encryption
- Between client and server for app interaction
- Between database servers for replication

TLS Encryption



Encryption All Storage

- Encryption at rest
- Integrate with external Key Management Systems
- Ability to rotate keys on demand

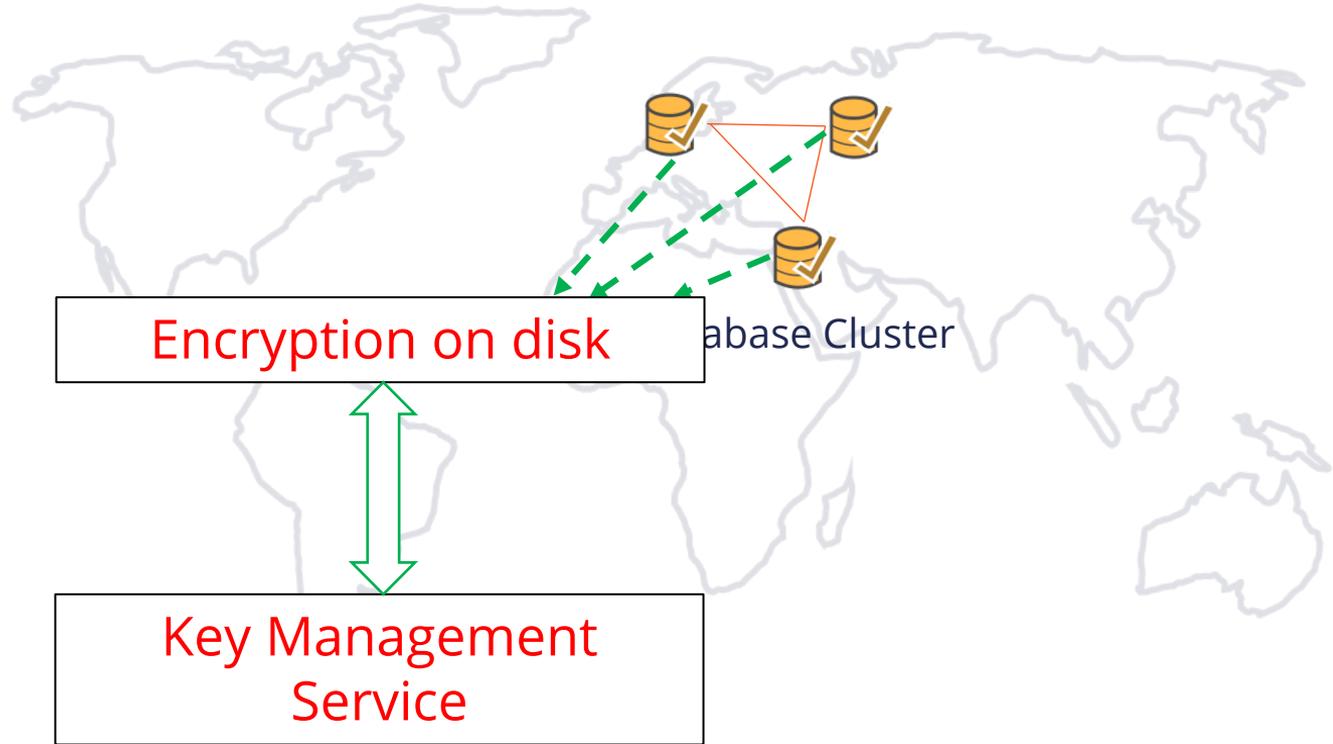


Use app level encryption if needed. Have a key-value table with id to cipher key. Encrypt PII data with the cipher key for fine-grained control. More in the next section.

Encryption at Rest



User



#3 RIGHT TO BE FORGOTTEN

Data must be erased on explicit request or when data is no longer relevant to original intent.



Use Encryption of Data Attributes

- Have a key-value table with id to cipher key
- Encrypt PII data with the cipher key on write
- Decrypt PII data on access
- Delete cipher key to forget PII data

Example - Storing User Profile Data

```
SET email=foo@bar.com FOR USER ID=XXX
```

Get encryption
key for user



- Reads require decryption
- Data not accessible without key

Encryption PII Data

```
SET email=ENCRYPTED FOR USER ID=XXX
```

```
CREATE TABLE users(  
  id          bigint PRIMARY KEY,  
  created_at text,  
  name       text,  
  email      text,  
  address    text,  
  city       text,  
  state      text,  
  zip        text,  
  birth_date text,
```

Store encrypted data

Use Anonymization of Data Attributes

- Many cases where value not needed
- Anonymize PII data with one way hash functions
- Use hashed ids for in data warehouse
- There is no PII data if hashed ids are used!

Example – Website Analytics

USER=foo@bar.com CHECKED OUT PRODUCT=X, CATEGORY=Gadget



One-way hash
user id

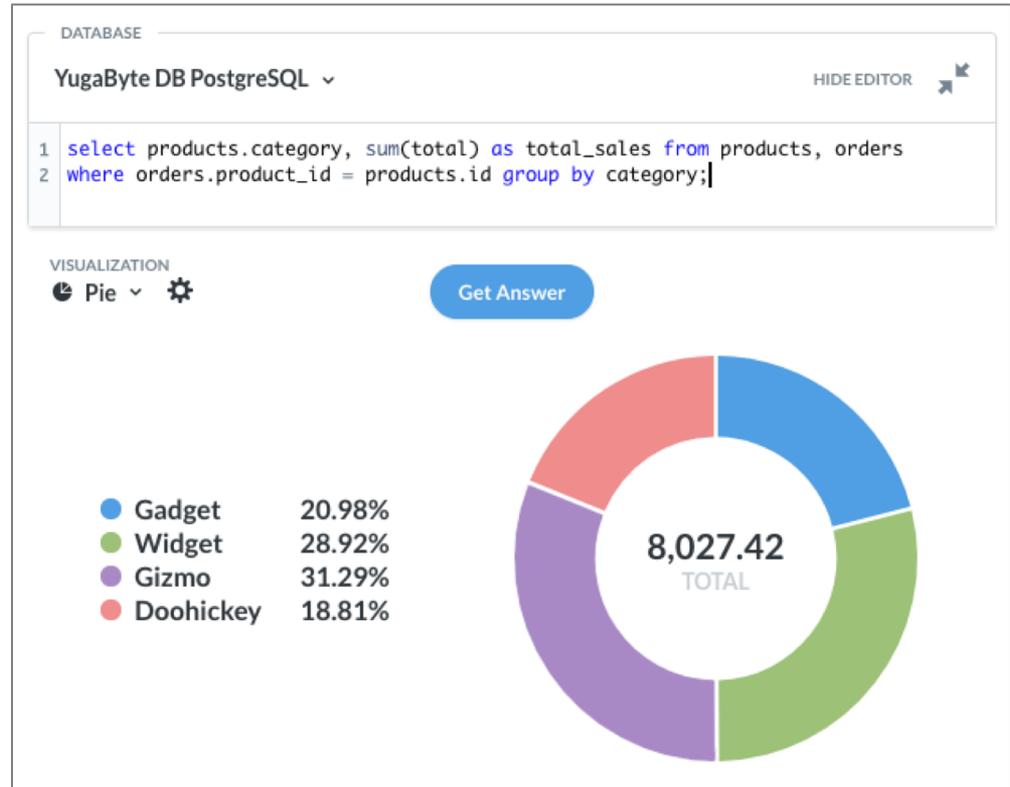


Analytics

USER=**HASHED_VAL** CHECKED OUT PRODUCT=X, CATEGORY=Gadget

Example – Website Analytics

- User no longer identifiable
- Hashed data still useful!



#4 DATA ACCESS ON DEMAND

Ability to inform a user about what data is being used, for what purpose and where it is stored.

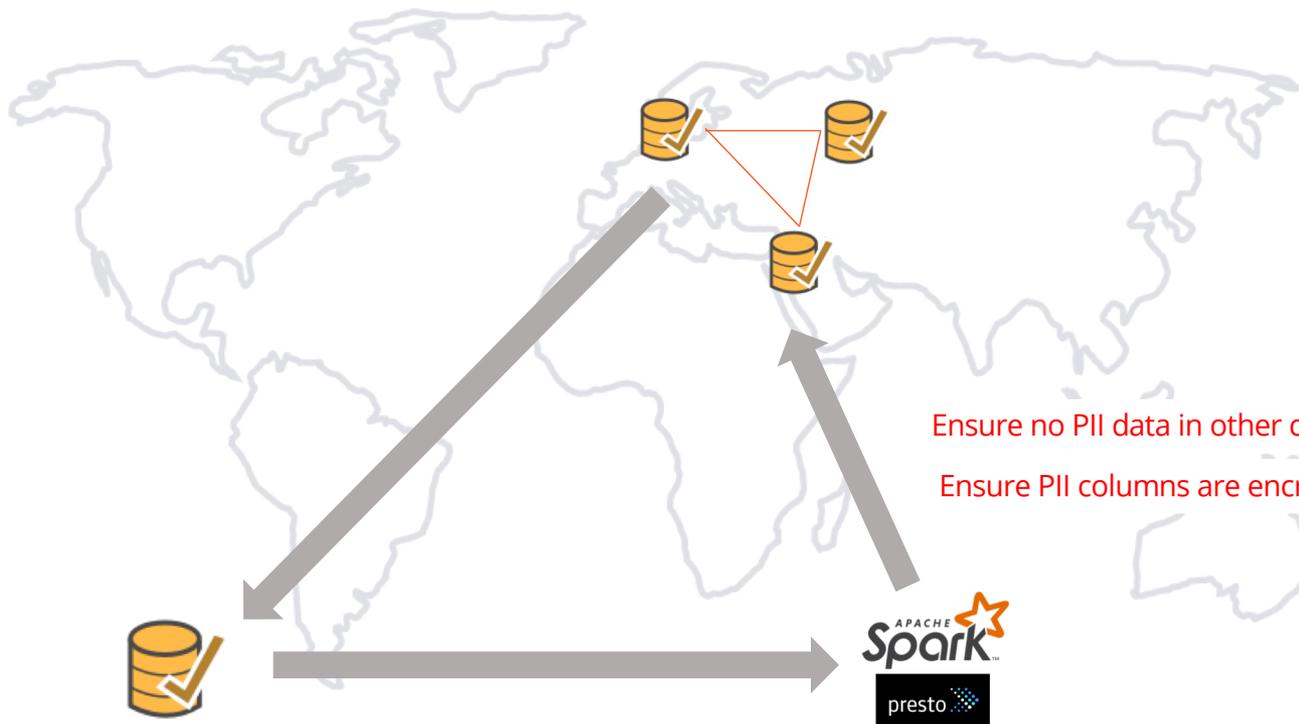


Tag Tables and Columns with PII

- Store in a separate information architecture table
- Make tagging a part of the process
- Easy to find what PII data is stored on demand

Run Continuous Compliance Checks

- Ensure PII are encrypted
- Ensure non-PII columns do not have sensitive data
- Use Spark/Presto to perform scan periodically
- Run scan on a read replica to not impact production

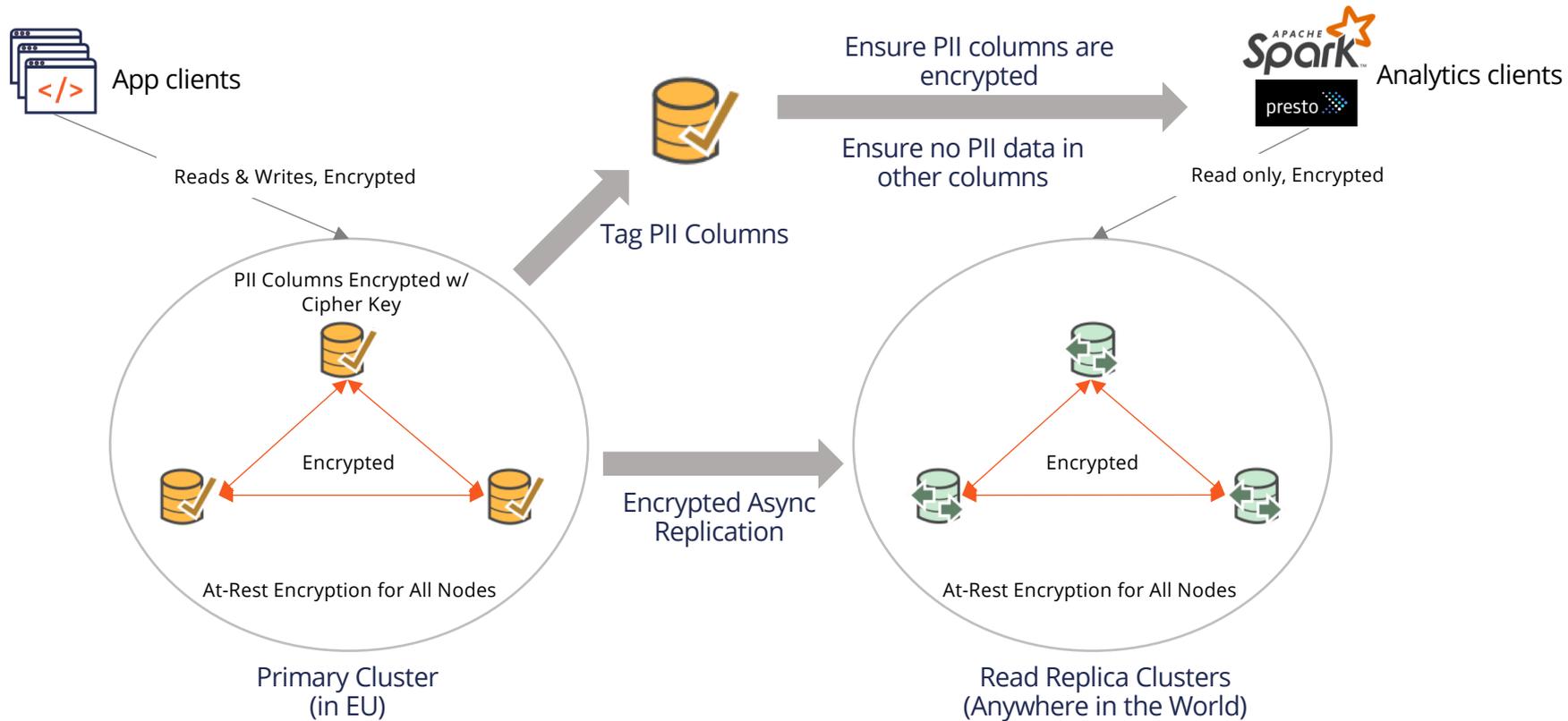


Tag PII Columns

Ensure no PII data in other columns
Ensure PII columns are encrypted

PUTTING IT ALL TOGETHER

GDPR Reference Architecture





Thank You!

Try it at

docs.yugabyte.com/latest/quick-start