# Advances in High Availability for PostgreSQL in the Enterprise

John Dalton
Senior Director, Product Management

August 25, 2021

EDB

# What you will learn today

1. Top-tier enterprise applications can run with confidence using Always On Postgres clusters

2. The latest Postgres-BDR™ release delivers faster throughput, larger clusters, and better observability

3. Applications requiring Oracle SQL compatibility and 99.999% availability can run on Postgres

# Agenda

- Enterprise needs for HA Postgres

- Faster throughput

- Larger clusters for data distribution

- Better observability

- HA for Oracle SQL compatible applications

- A look ahead

# What is "Always On?"

Delivering mission-critical applications/services 24x7

## Finance

Payments

Bank account access

## Telecom

Video conferencing

Texting/alerting

## Transportation

Travel

Rideshare

GPS

# Do you need to be "Always On"?

## Things to consider

What's the reputation cost of downtime to your business?

Are there times when it's okay for an application to be inaccessible?
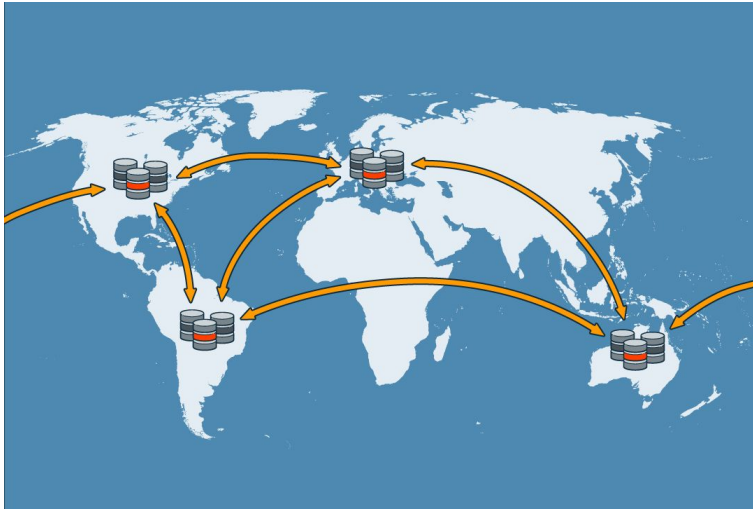
Is access to data tied directly to revenue?

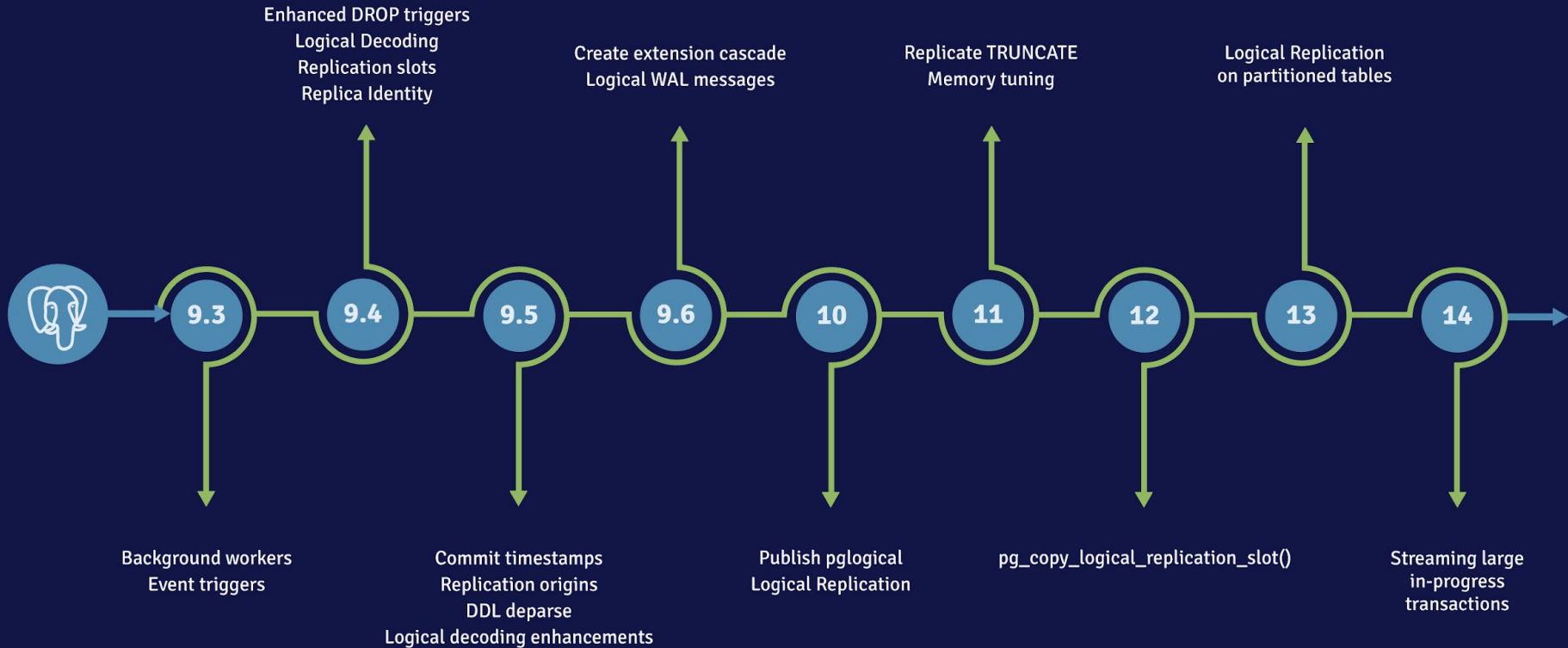Are your customers globally distributed?

# Postgres-BDR is more than bi-directional replication

**Multi-master replication enabling highly available and geographically distributed Postgres clusters**



- Logical replication of data and schema enabled via standard PostgreSQL extension

- Data consistency options that span from immediate to eventual consistency

- Robust tooling to manage conflicts, monitor performance, and validate consistency

- Deploy natively to cloud, virtual, or bare metal environments

Postgres **BDR**

Enhanced DROP triggers
Logical Decoding
Replication slots
Replica Identity

Create extension cascade
Logical WAL messages

Replicate TRUNCATE
Memory tuning

Logical Replication
on partitioned tables

| 9.3 | 9.4 | 9.5 | 9.6 | 10 | 11 | 12 | 13 | 14 |

Background workers
Event triggers

Commit timestamps
Replication origins
DDL deparse
Logical decoding enhancements

Publish pglogical
Logical Replication

pg_copy_logical_replication_slot()

Streaming large
in-progress
transactions

# That's great, but enterprises are asking for more

And maintaining data consistency is paramount

## Faster Throughput

I need to push more transactions (TPS) through the cluster, as fast as Postgres can run

## Larger Clusters

My geographically distributed application needs a read scalable cluster without increased overhead

## Better Observability

How do I ensure the right data is exposed to quickly triage issues and facilitate recovery actions

## More Migrations

I want to move my applications that need 99.999% availability from Oracle to Postgres

# Faster Throughput

# Overview

Introducing parallel data flow architecture to BDR

Goals of pipeline parallelism

- Smoother latency
- Less CPU overhead on BDR upstream nodes
- Better overall throughput for parallelizable workloads



## Single decoding worker [upstream]

Single decoding worker improves performance on upstream node by doing logical decoding of WAL once instead of for each downstream node
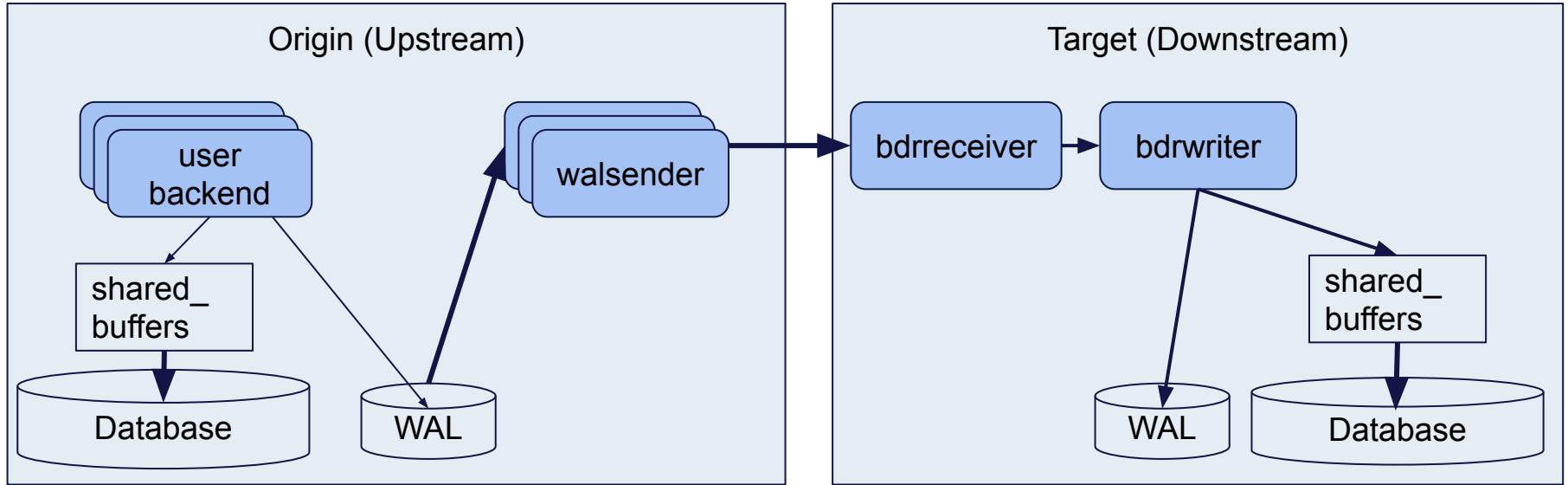
Currently not enabled by default, expected to quickly mature

## Parallel apply [downstream]

Parallel apply allows multiple writer processes to apply transactions on downstream node with throughput up to 5X faster
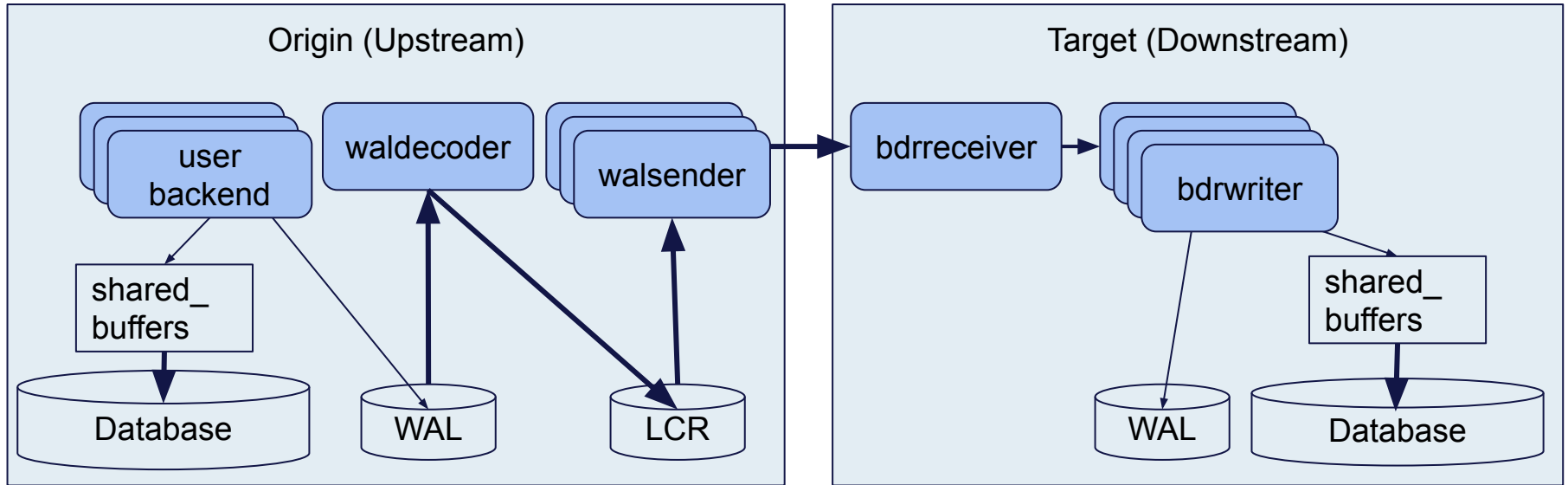
Benefits mixed write workloads, especially for larger DBs and/or i/o heavy write workloads

# Logical streaming replication in 3.6

# Parallel logical streaming replication in 3.7
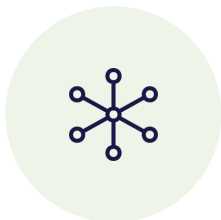
## Only available with BDR

# Larger Clusters

# Introducing data distribution

Support for a new use case enabling larger clusters

## High read scalability is required

For geographically distributed applications requiring consistent, shared reference data

For example Telecom call routing, this can mean **hundreds of nodes** in a cluster across different regions
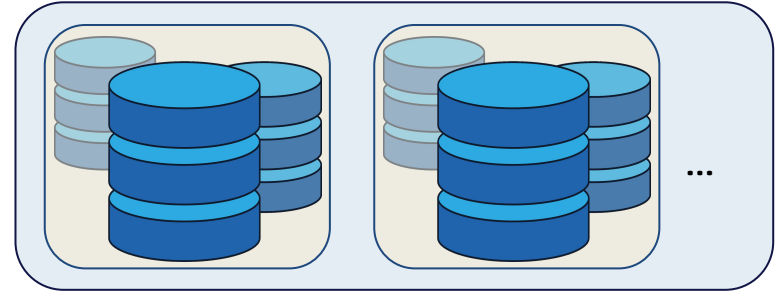
## Built on new enabling capabilities

New features supporting this use case are sub-groups and subscribe-only nodes

Also leveraged is single decoding worker to minimize overhead on origin node

# Sub-groups
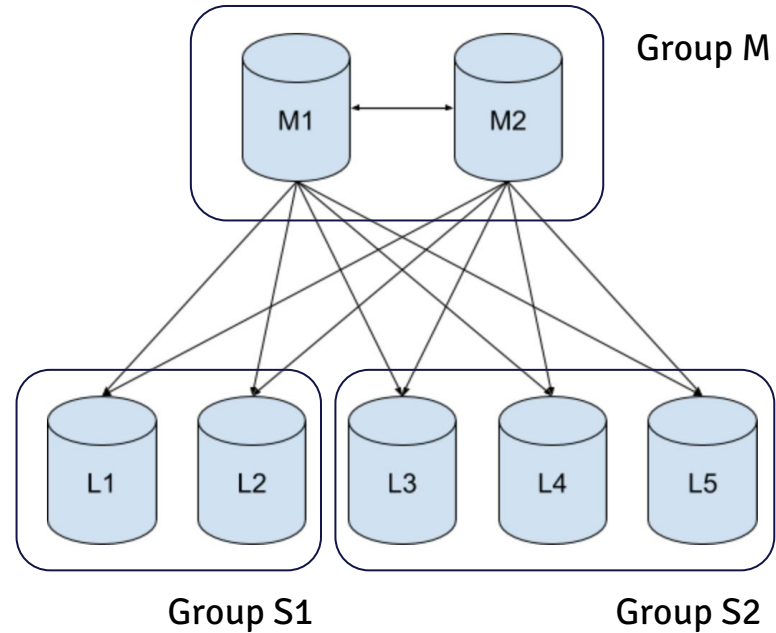
## Foundational feature for many planned enhancements

- In BDR3.6, all nodes were part of one Group
  - Mesh architecture minimizes latency between nodes

- In BDR3.7, a cluster can have nested SubGroups, with each node as part of one Group
  - AlwaysOn becomes 2 Group cluster (2x2)
  - Each Group can have >=2 Nodes
  - No limit on number of Groups
  - Each Group comes with own default repset to simplify management of filtering

- **This is a building block for more complex architectures**

# Data distribution tree

## Enables sharing reference data for applications requiring high read scalability

- Tree architecture minimizes **overheads in large networks**

- Data Distribution Networks

- Subscriber-Only Nodes only **receive** data from Main Group(s)

- Subscriber-Only Nodes in separate groups (sub-groups)

- Allows up to **1000** node clusters

Group M

M1 ↔ M2

L1  L2    L3  L4  L5

Group S1            Group S2

# Better Observability

# Monitoring BDR

Postgres Enterprise Manager (PEM) now includes visualizations for BDR

## PEM dashboards

The latest release of Postgres Enterprise Manager introduces 3 BDR monitoring dashboards for displaying information about replication activity for a BDR cluster. These are:

- Admin
- Group Monitoring
- Node Monitoring

## BDR enhancements

Ongoing work to improve operational insights of deployed BDR clusters. With latest BDR release a number of new views are introduced with focus on:

- Group level monitoring
- In-progress monitoring on the downstream apply side

# BDR admin dashboard

This view provides cluster (or group) level information on:

- Node status across the cluster

- Status on global locks (DDL)

- Software version details by node

- Apply side worker thread status

- Worker thread errors

- CAMO status for cluster

- Raft consensus status for cluster

# BDR group monitoring dashboard

Insight to overall cluster (or group) operational health of replication is available in this view. Regular advancement of group level metrics means all nodes are actively consuming changes. Information available includes:
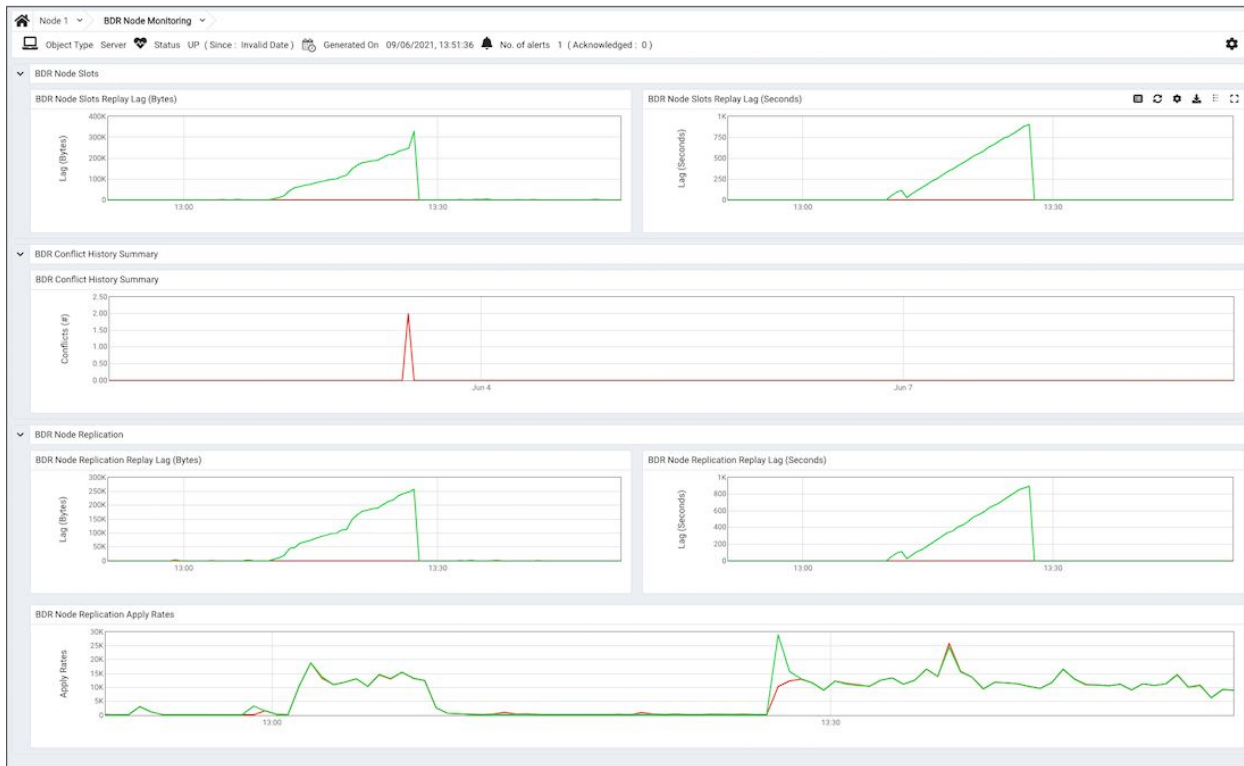
- Subscription lag

- Replication slots replay, flush, and write lag in bytes and seconds

- Replication slots sent lag in bytes

# BDR node monitoring dashboard

Node level visibility:

- For slots used by BDR in a database the outgoing replication replay lag in bytes and seconds (difference between applied LSN and current WAL write position)

- A summary of conflicts encountered

- Outgoing replication rates for a node including replay lag in bytes and seconds

- Rate in LSNs applied per second for a BDR node

# HA for Oracle SQL Compatible Applications

# EDB Postgres Advanced with BDR

Benefits of Oracle SQL compatibility and 5 nines availability

## Your tier 1 apps

Leverage existing infrastructure with native PL/SQL support and Oracle Call Interface (OCI) interoperability

## Your people

Leverage the existing skills of your Oracle DBAs and developers

## Your business

Lower costs, reduce risks, and move faster

# BDR Feature Overview

A full-featured multi-master replication solution for PostgreSQL clusters

## Essentials

**Provides the essential multi-master capabilities for PostgreSQL clusters.**

- Enables application and database upgrades without requiring downtime

- Provides clusters row level eventual consistency by default

- Tools to monitor operation and verify data consistency

- Extends PostgreSQL logical replication beyond unidirectional, standby use cases

## Advanced

**Includes advanced conflict management, data-loss protection, and up to 5X faster.**

- Guards applications from committing transactions more than once

- Conflict-free synchronous replication with two phase commit

- Concurrent updates using conflict-free replicated data types (CRDTs)

- Configurable column level conflict resolution with customizable conflict handling and transformation

# A Look Ahead

# Roadmap

A vision for advancing very highly available Postgres clusters

## Q2: BDR 3.7 GA

Support for v11-13

Up to 5X throughput with parallel apply

Single decoding worker streamlines upstream replication

High read scalability through data distribution tree architecture with subscribe-only nodes

Oracle SQL compatibility support with Postgres Advanced

PEM support for monitoring BDR

## Q4: BDR 4.0 GA

Support for v12-14

Oracle SQL compatibility support is complete with Postgres Advanced v14

Features enabled are:

- Eager replication
- CAMO
- Single decoding worker
- Application assessment

Improved cluster management provides zero transaction lag switchover

## 2022 and beyond

Support for v13-15

Expand pipeline parallelism and enhance large transaction support for faster throughput

Enhance autopartition to support bigger tables in clusters

Evolve data distribution capabilities in support of bigger clusters

Develop autoscale capabilities to enable bigger databases

# Thank you!

john.dalton@enteprisedb.com

EDB™