

The background features a dark navy blue field with abstract, overlapping shapes in vibrant magenta and deep red. Two thin, light blue lines intersect diagonally across the upper right portion of the image. The text is positioned on the left side.

AWS re:Invent

DECEMBER 2 – 6, 2024 | LAS VEGAS, NV

DAT420

Achieving scale with Amazon Aurora PostgreSQL Limitless Database

Anum Jang Sher

Senior Product Manager
Amazon Aurora

David Wein

Senior Principal Technologist
Amazon Aurora



© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Agenda

01 Scaling challenges

02 Overview

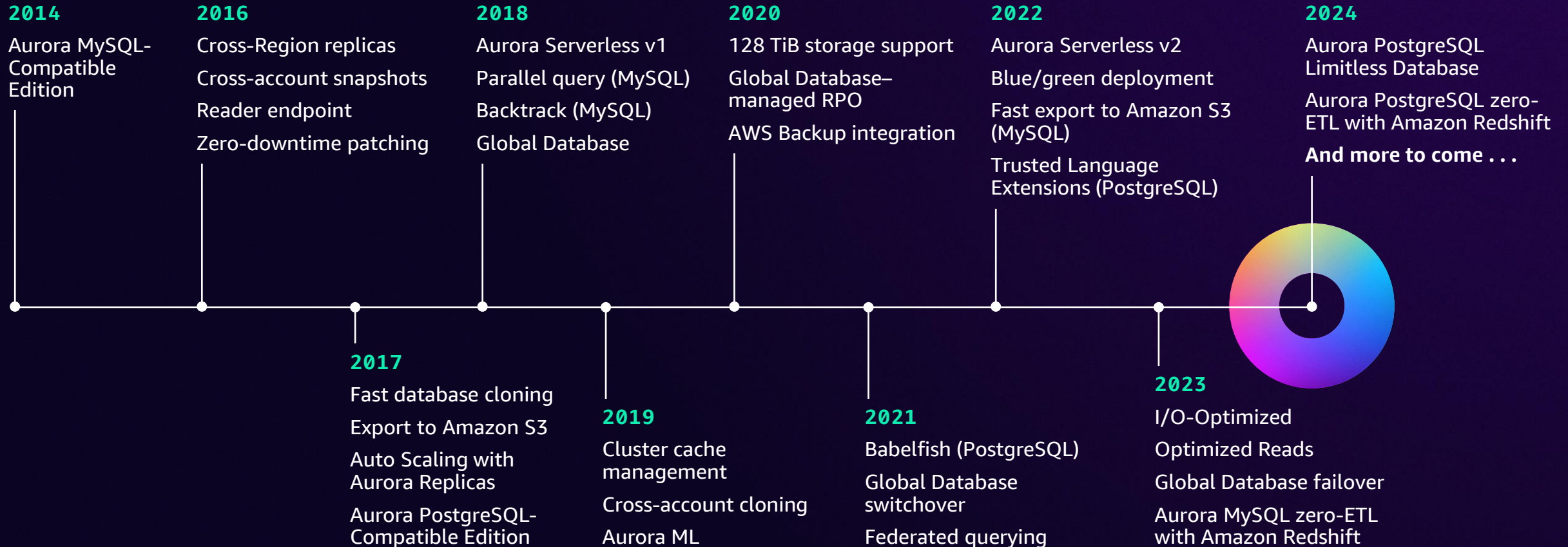
03 Architecture

04 Data distribution

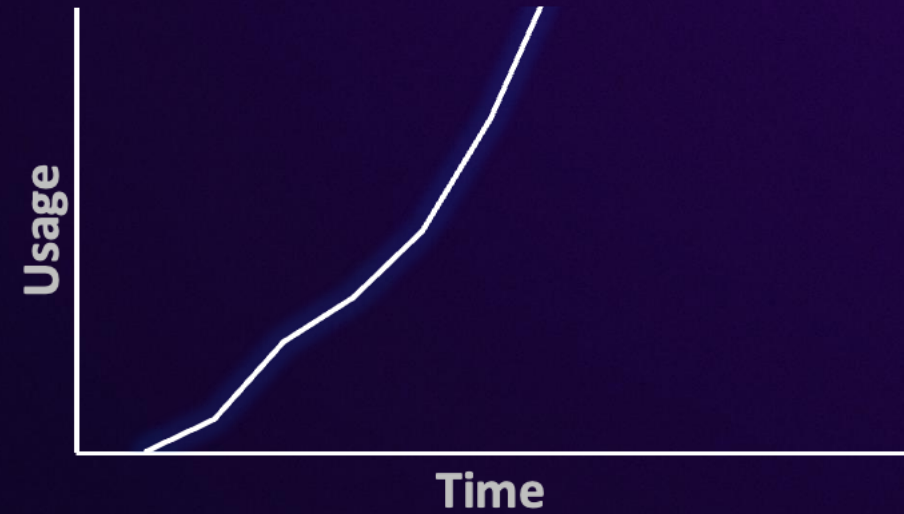
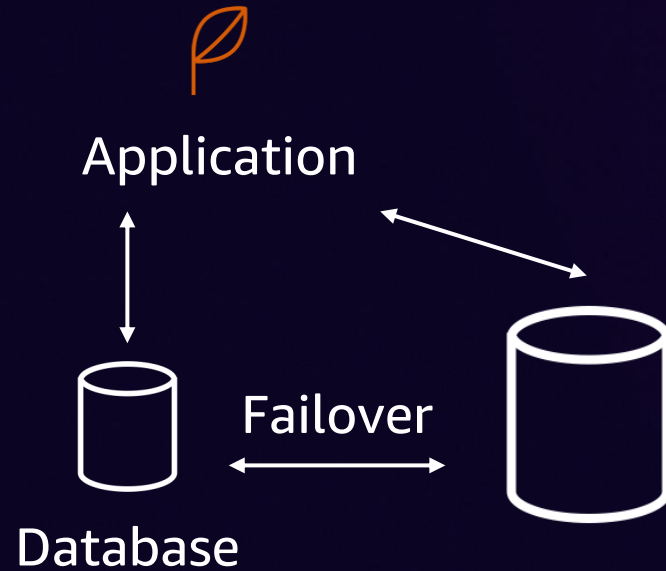
05 Transactions and queries

06 Get started today!

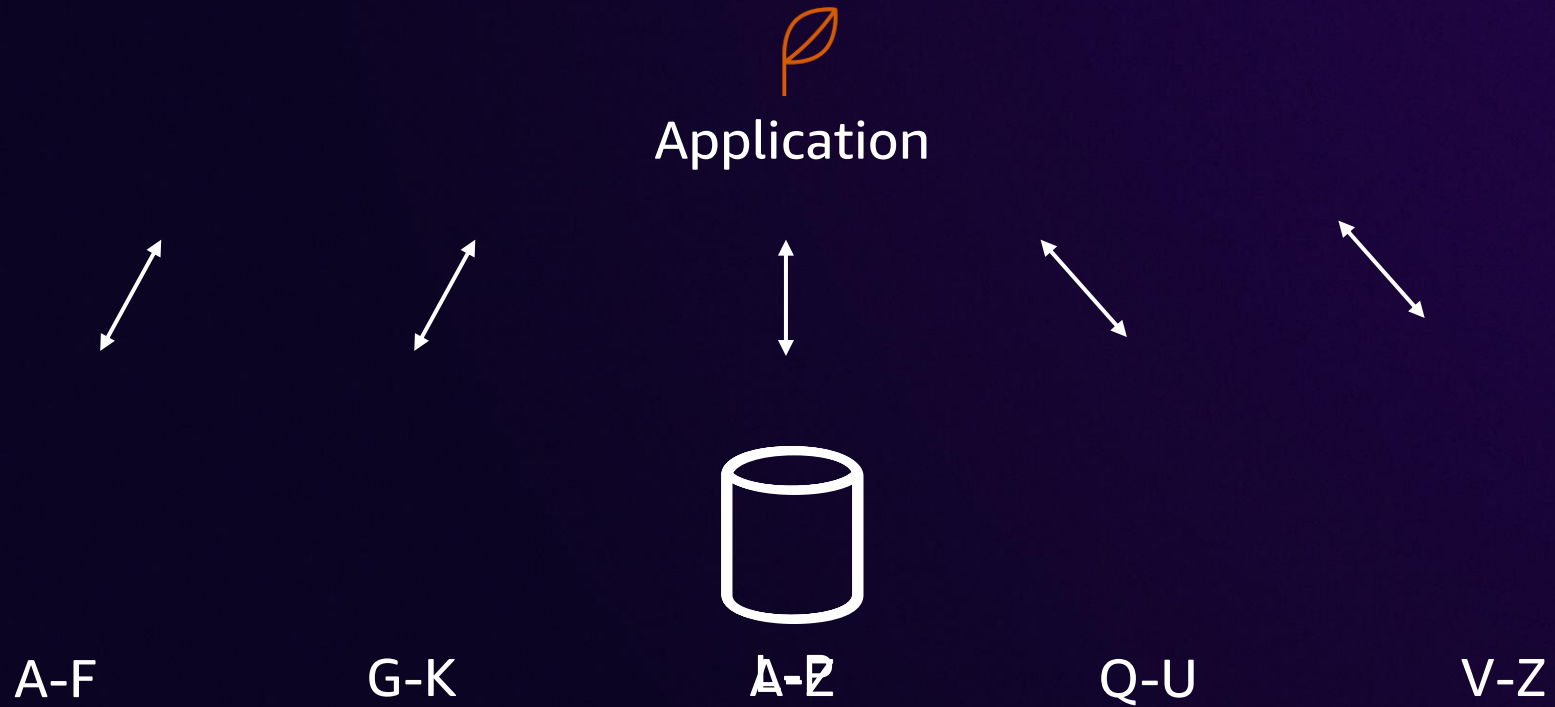
Celebrating a decade of Amazon Aurora innovation for customers



Scaling databases



Sharding brings **scale**



Challenges



Querying



Consistency

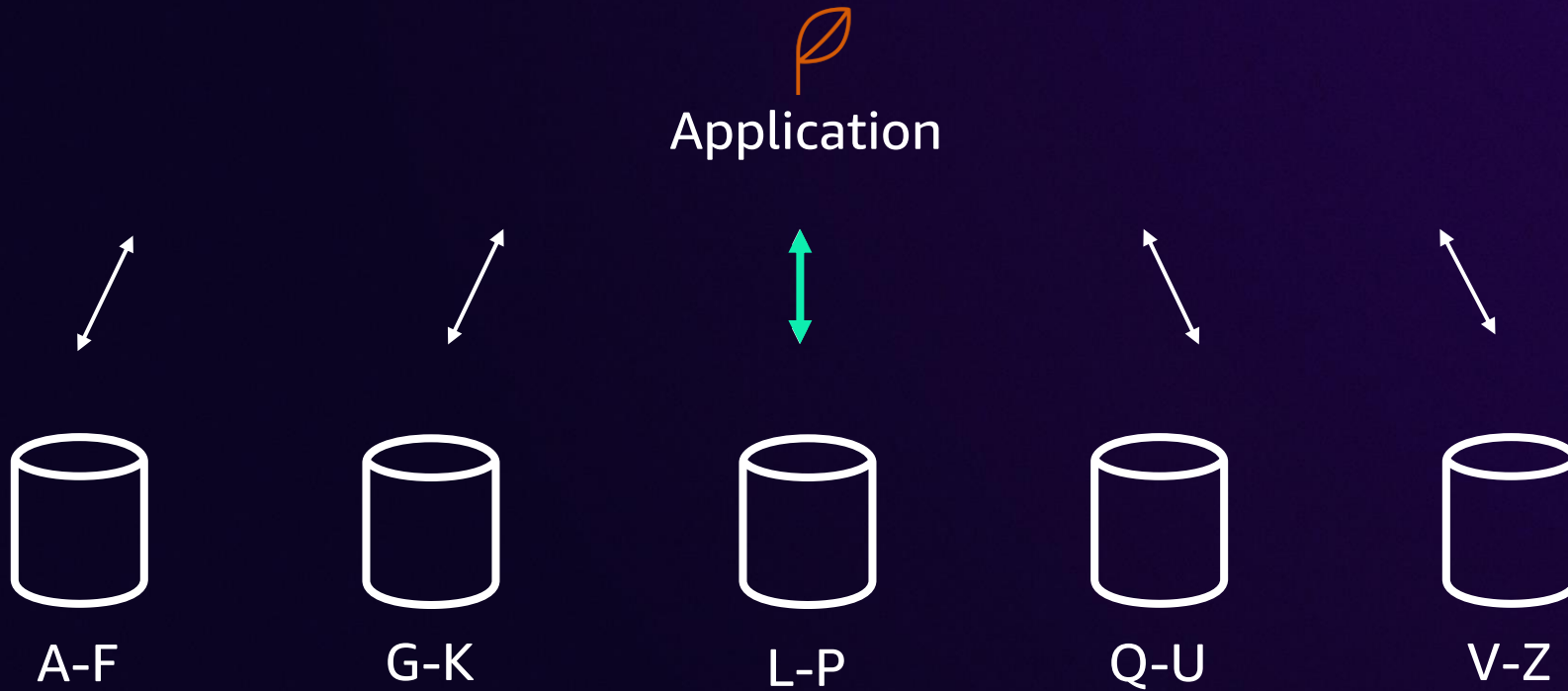


Re-sharding

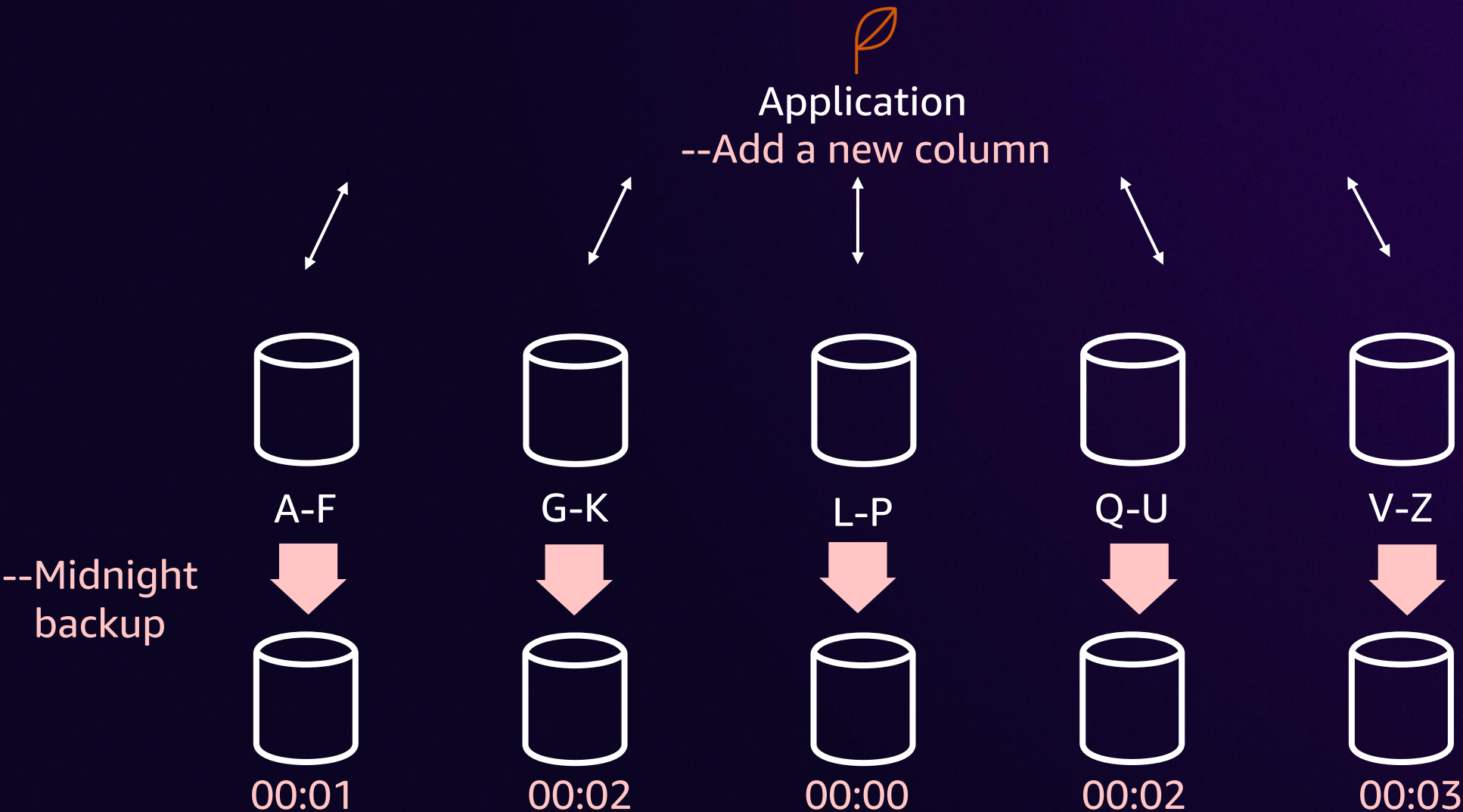


Capacity management

Challenges: Querying

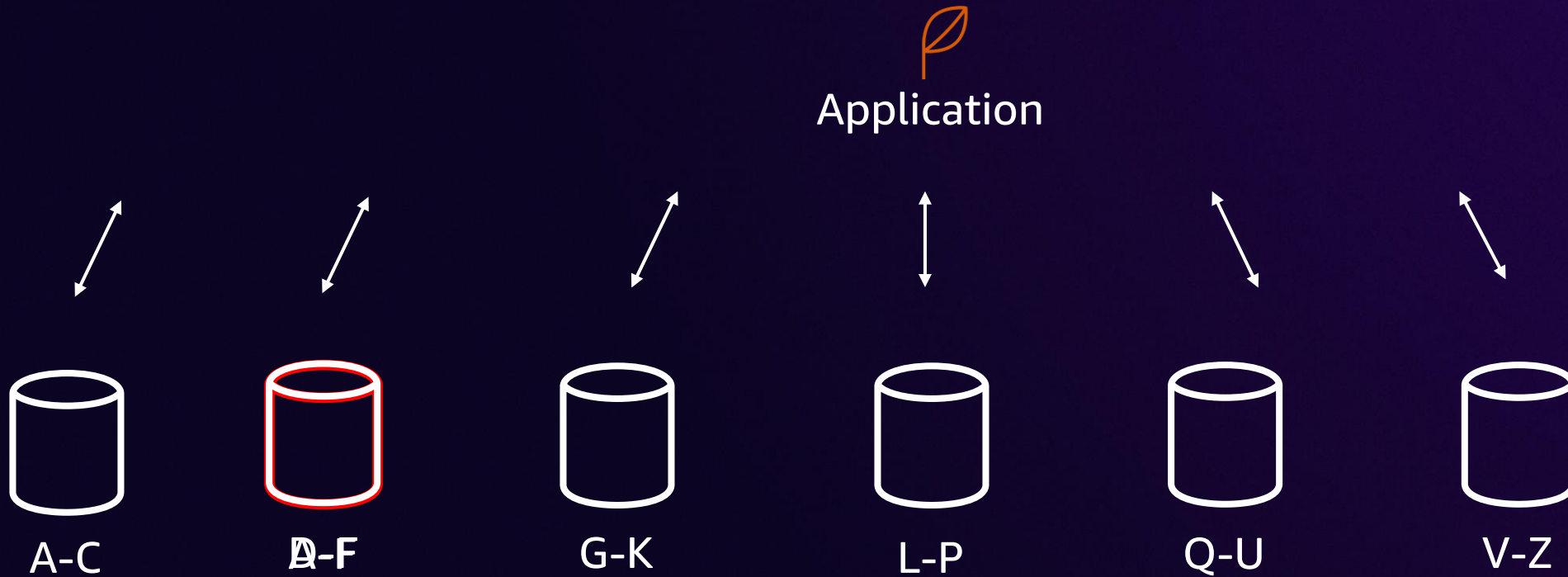


Challenges: Consistency

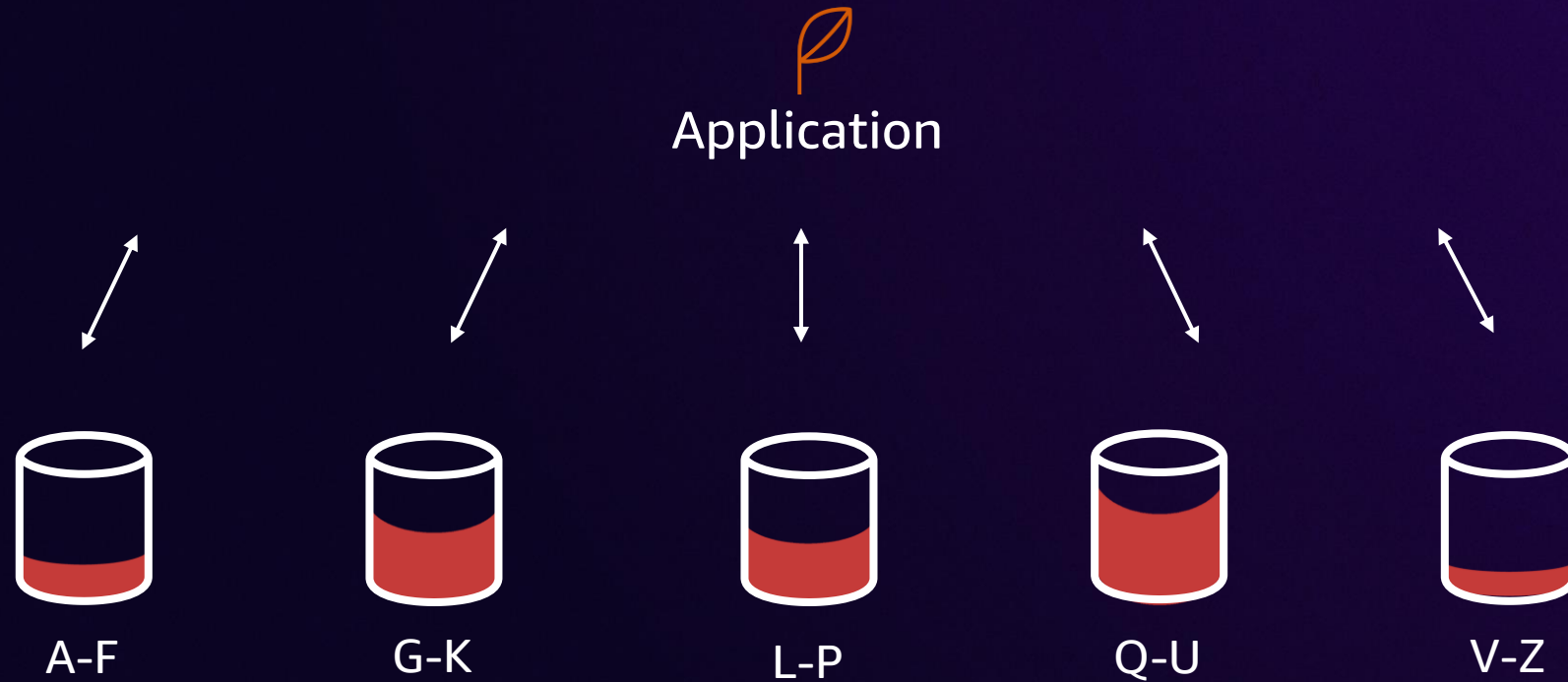




Challenges: Re-sharding



Challenges: Database capacity management



Aurora PostgreSQL Limitless Database

MANAGED HORIZONTAL SCALE-OUT BEYOND THE LIMITS OF A SINGLE INSTANCE



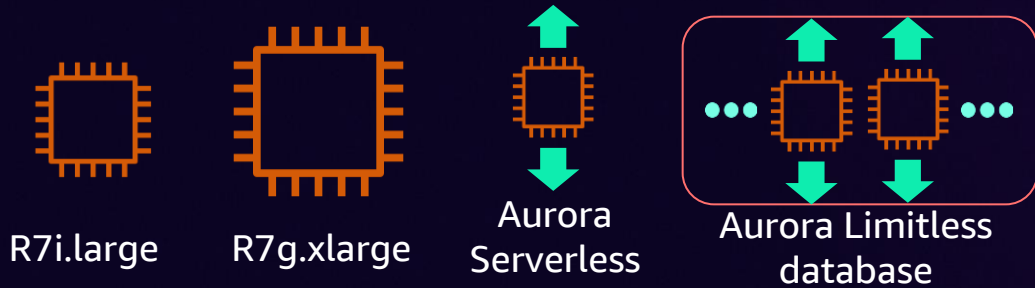
Scales to **millions of write transactions** per seconds

Manages **petabytes** of data

Operational simplicity of a **single** instance

Pay-per-use pricing

DB shard group



Aurora Limitless Database - new [Info](#)

With Limitless Database, Aurora can automatically scale write throughput and data storage capacity beyond the limits of a single DB cluster.

DB shard group identifier
Type a name for your DB shard group. The name must be unique across all DB shard groups owned by your AWS account in the current AWS Region.

Constraints: 1 to 60 alphanumeric characters or hyphens. First character must be a letter. Can't contain two consecutive hyphens. Can't end with a hyphen.

DB shard group capacity range [Info](#)
Enter the minimum and maximum capacity for Limitless Database. The capacity is measured in Aurora capacity units (ACUs) across all routers and shards.

Minimum capacity (ACUs) (48 GiB) **Maximum capacity (ACUs)** (1040 GiB)

Enter a value greater than or equal to 16 ACUs Enter a value less than or equal to 6144 ACUs

DB shard group deployment
The number of additional cross Availability Zone standby shards. Adding compute redundancy will have a significant impact on cost. [Learn more](#)

☐ No compute redundancy
Creates a DB shard group without standbys for each shard.

☒ Compute redundancy with a single failover target
Each shard is created with one compute standby in a different Availability Zone.

☐ Compute redundancy with two failover targets
Each shard is created with two compute standbys in different Availability Zones.

Create **DB shard group** instead of instances

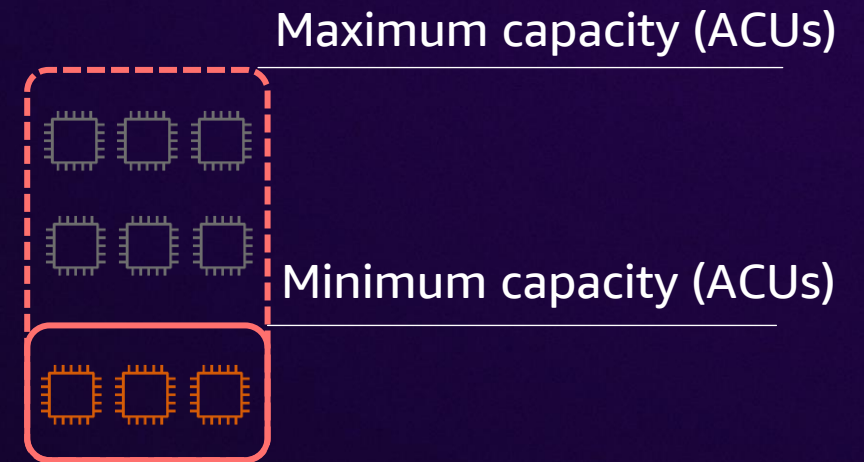
Automatic scaling based on workload

Specify **compute redundancy**

Supported by **99.99%** availability SLA

Capacity management

- DB shard group capacity is measured in **Aurora Capacity Units** (ACUs)
- One ACU is a combination of **2GiB** of memory
- Corresponding **CPU** and **networking**
- You set the **minimum** and **maximum** capacity



Aurora PostgreSQL
Limitless Database

Scenario

Customer

<input type="checkbox"/>	<input type="text"/>
<input type="checkbox"/>	<input type="text"/>
<input type="checkbox"/>	<input type="text"/>
<input type="checkbox"/>	<input type="text"/>

cust_id
name
email

Order

<input type="checkbox"/>	<input type="text"/>
<input type="checkbox"/>	<input type="text"/>
<input type="checkbox"/>	<input type="text"/>
<input type="checkbox"/>	<input type="text"/>

order_id
cust_id
amount
tax_rate_id

Tax rate

<input type="checkbox"/>	<input type="text"/>
<input type="checkbox"/>	<input type="text"/>
<input type="checkbox"/>	<input type="text"/>
<input type="checkbox"/>	<input type="text"/>

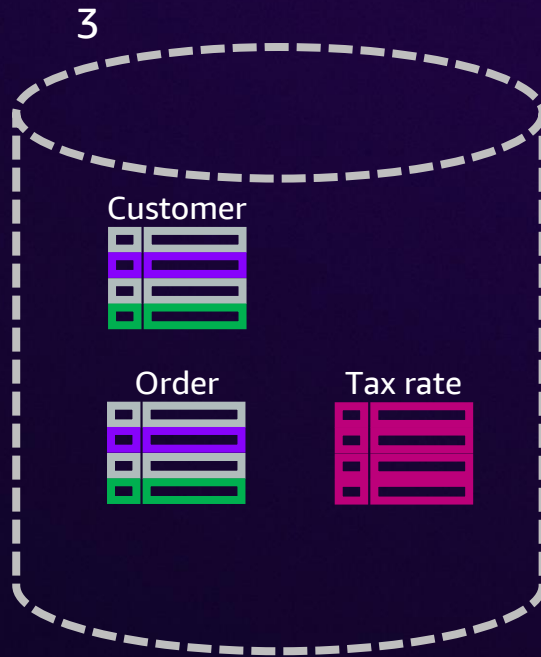
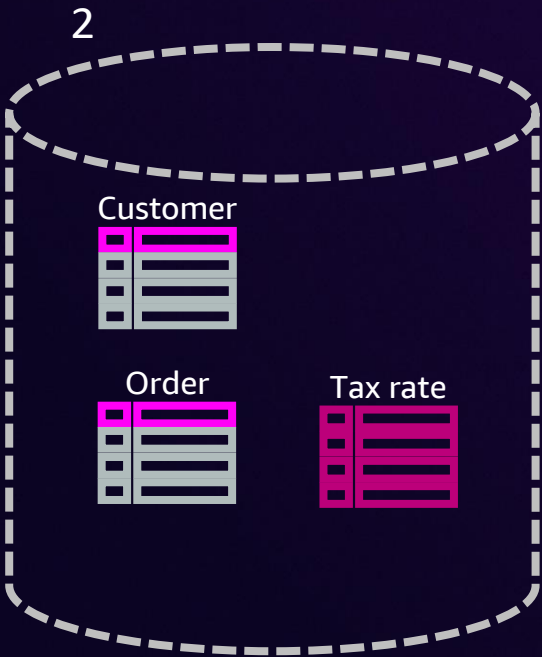
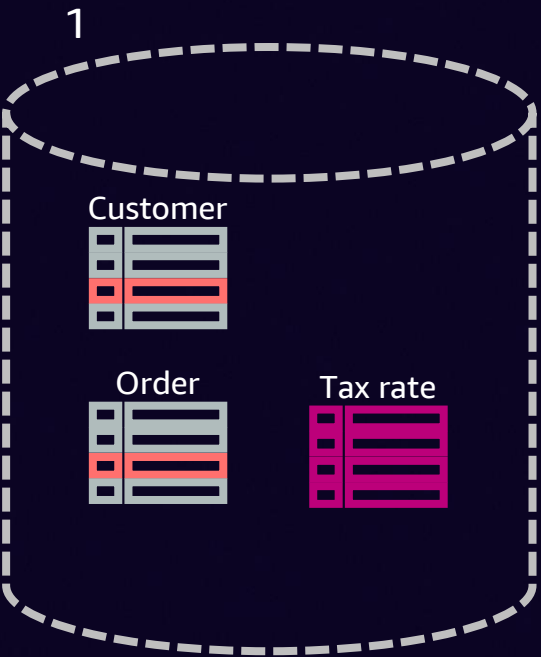
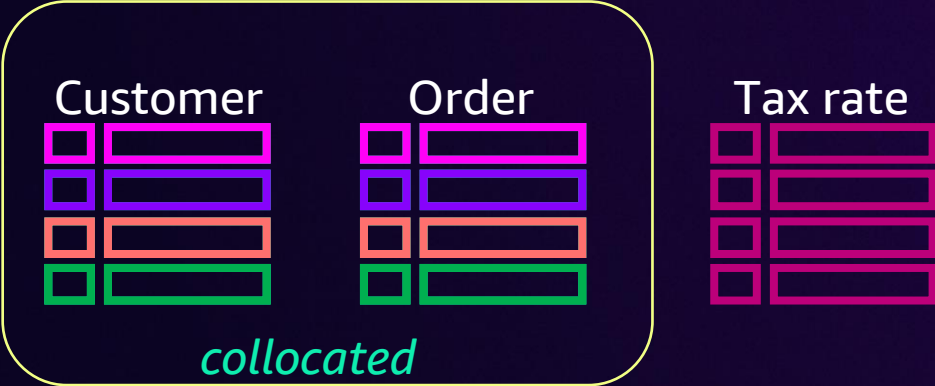
tax_rate_id
city
state
country
tax_rate

Use limitless database

Standard

Sharded

Reference



Create sharded customer table

```
SET rds_aurora.limitless_create_table_mode='sharded';
```

```
SET rds_aurora.limitless_create_table_shard_key='{“cust_id”}';
```

```
CREATE TABLE customer (  
    cust_id INT PRIMARY KEY NOT NULL,  
    name TEXT,  
    email VARCHAR(100)  
);
```

Create collocated order table

```
SET rds_aurora.limitless_create_table_mode='sharded';  
SET rds_aurora.limitless_create_table_shard_key='{“cust_id”}';
```

```
SET rds_aurora.limitless_create_table_collocate_with='customer';
```

```
CREATE TABLE order (  
    order_id    INT NOT NULL,  
    cust_id     INT NOT NULL,  
    amount      DOUBLE NOT NULL,  
    tax_rate_id DOUBLE,  
    PRIMARY KEY (order_id, cust_id)  
);
```

Create reference table tax_rate

```
SET rds_aurora.limitless_create_table_mode='reference';
```

```
CREATE TABLE tax_rate (  
    tax_rate_id INT PRIMARY KEY NOT NULL,  
    city        TEXT NOT NULL,  
    state       TEXT,  
    country     TEXT NOT NULL,  
    tax_rate    DOUBLE NOT NULL  
);
```


Amazon Aurora PostgreSQL Limitless Database architecture

David Wein

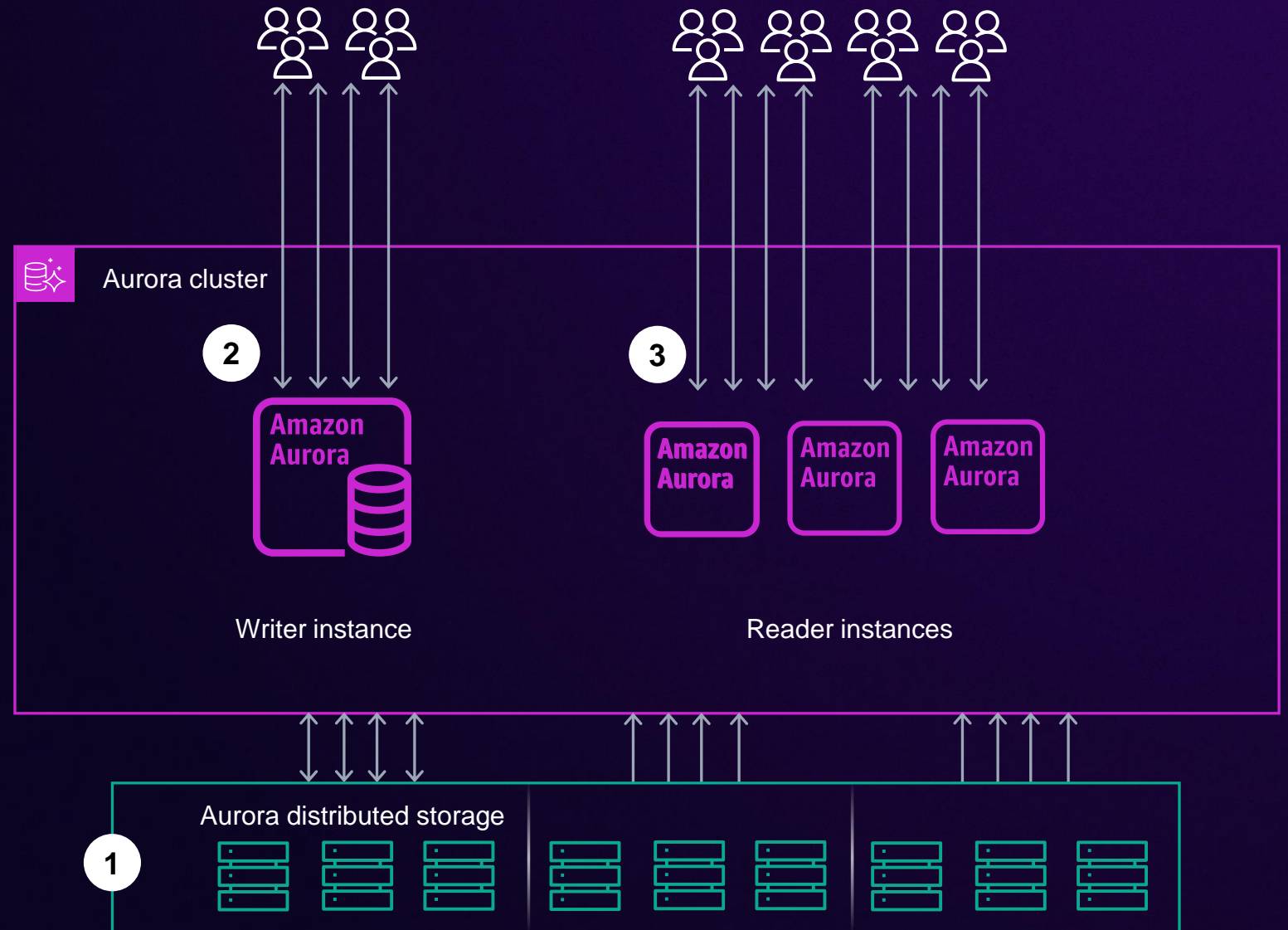
Senior Principal Technologist
Amazon Aurora



© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Standard Aurora architecture

- 1 Aurora volume on **distributed storage**
- 2 An Aurora **writer** instance
- 3 Optional **reader instances** for availability and read scaling



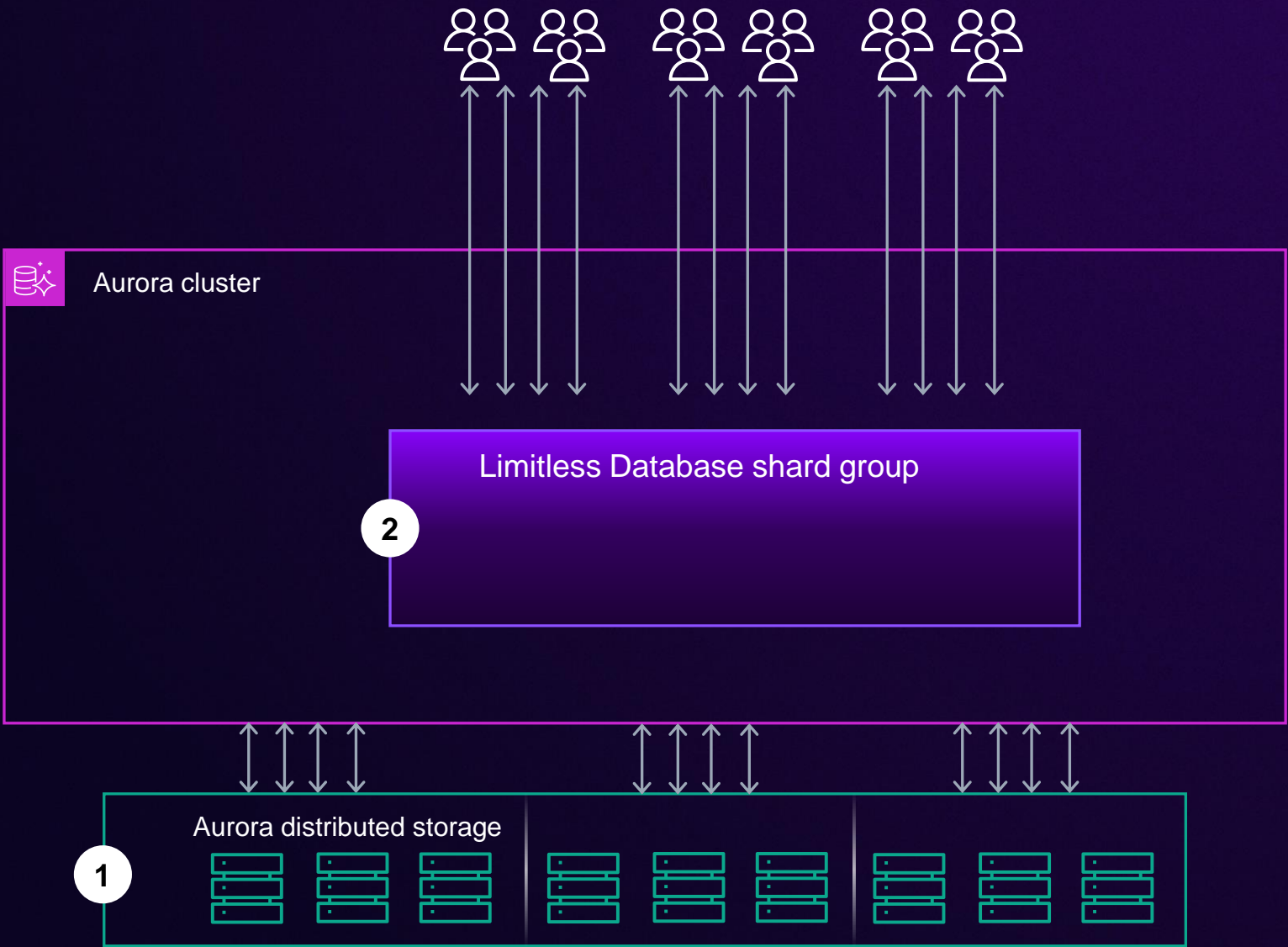
Standard Aurora architecture

1 Aurora volume on distributed storage

An Aurora **writer** instance

Optional **reader** instances for availability and read scaling

2 Limitless Database introduces the “**shard group**” concept



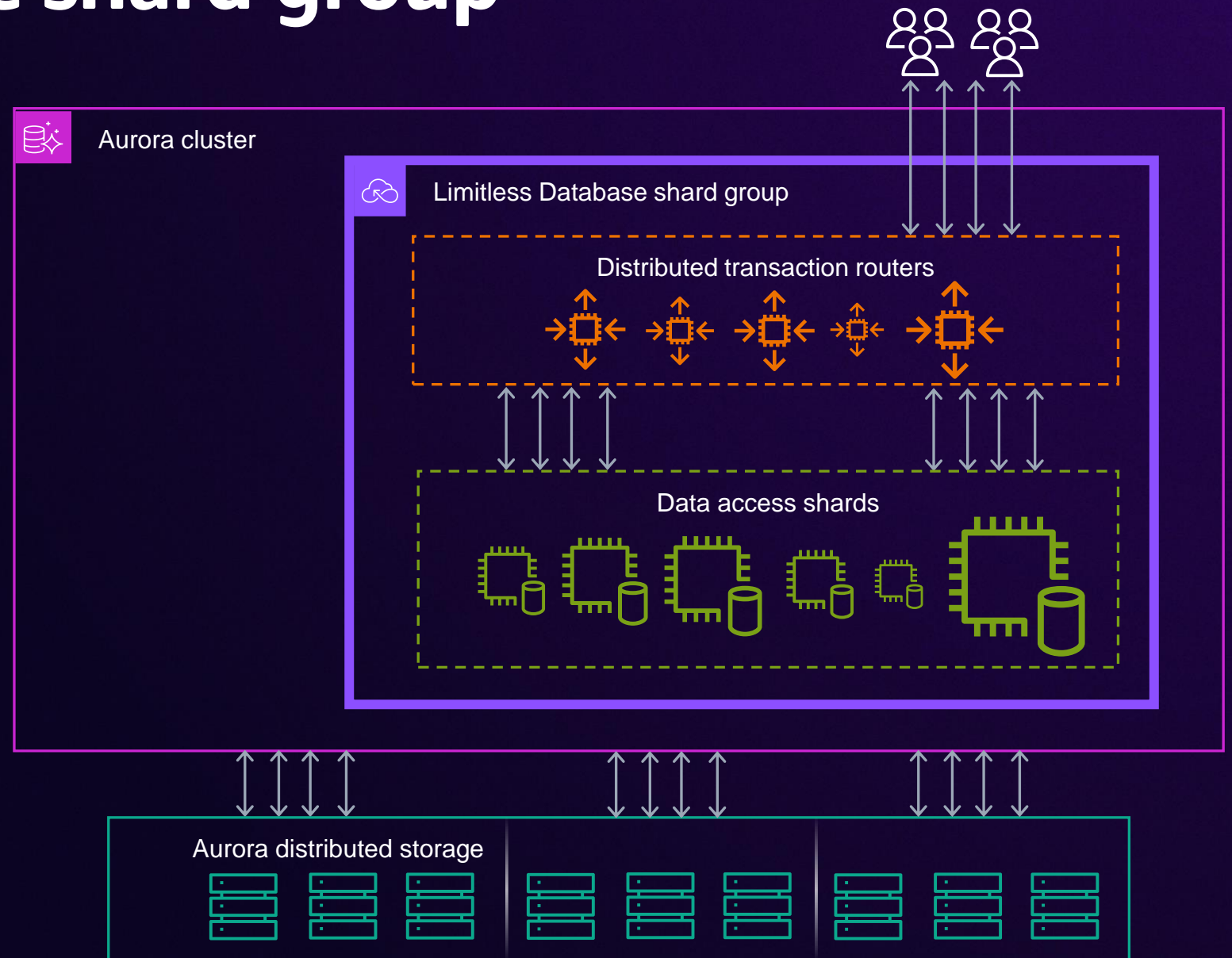
Limitless Database shard group

Contained within your **Aurora cluster**

Encapsulates limitless database infrastructure for your cluster

Provides **an endpoint** for applications

Scales resources within configured capacity based on **load** and **data size**



Distributed transaction routers

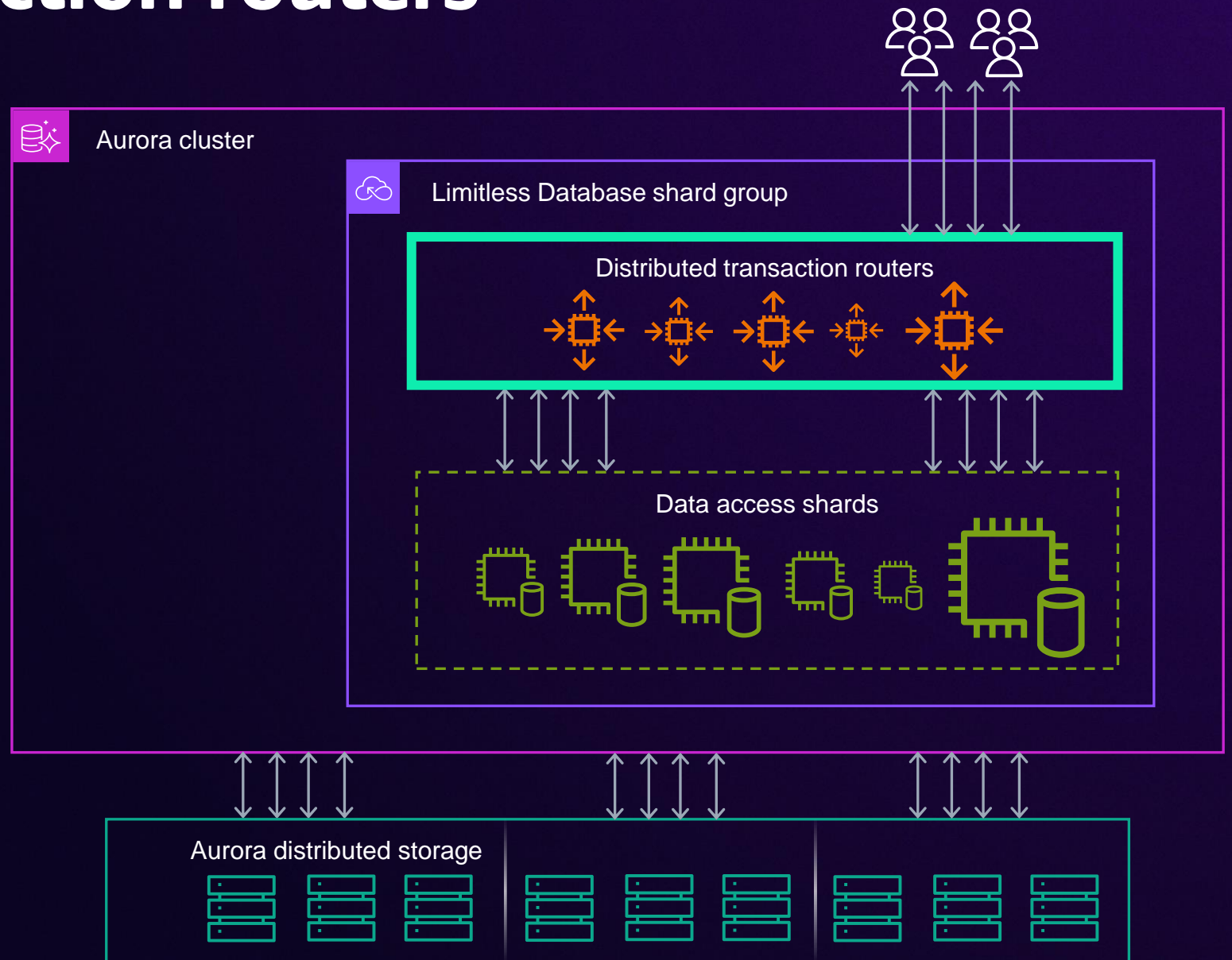
Serve **all** application traffic

Scale **vertically** and **horizontally** based on load

Know schema and key range placement

Assign time for transaction snapshot and drive **distributed commits**

Perform initial planning of query and aggregate results from **multi-shard queries**



Data access shards

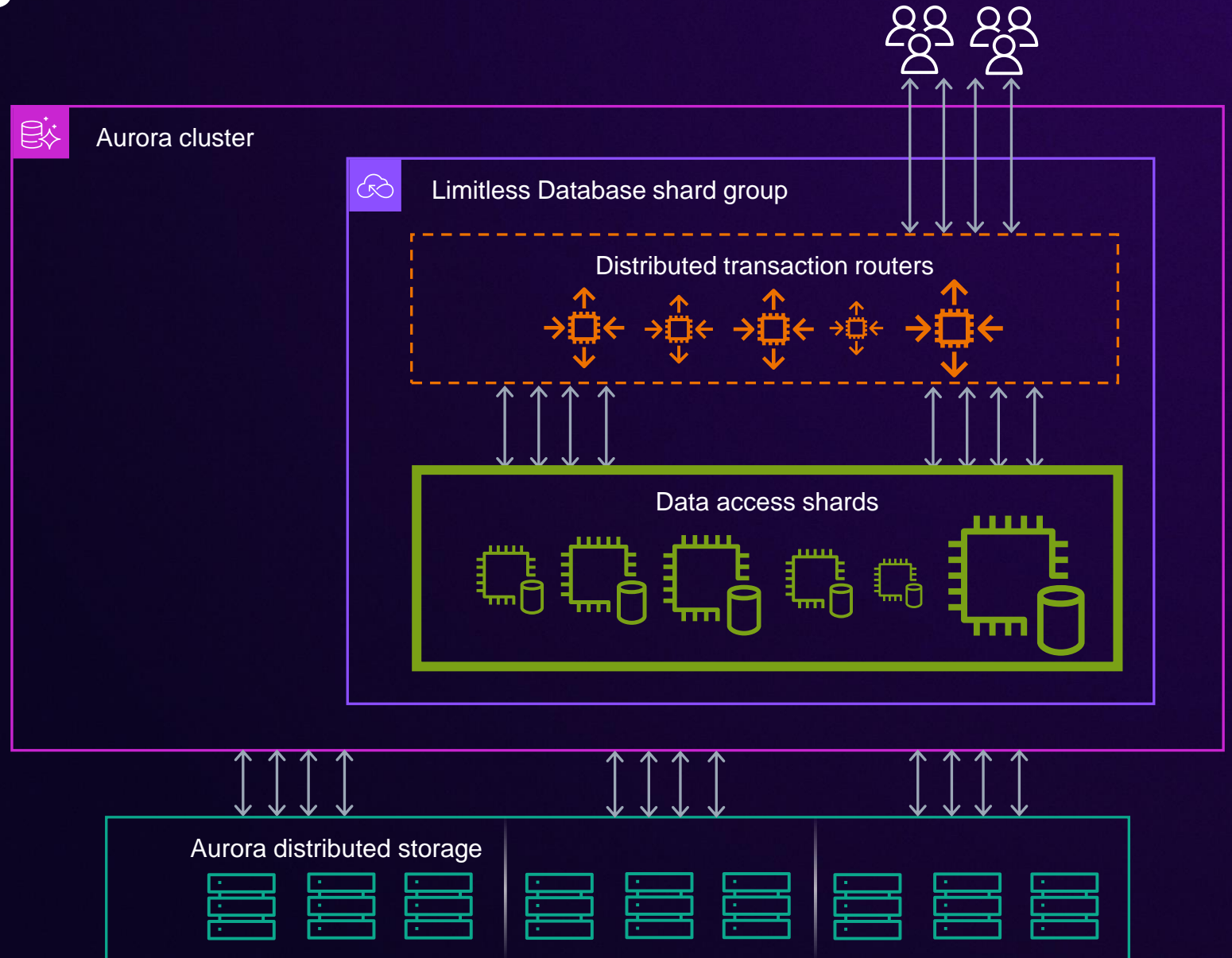
Own **portion** of sharded table key space and have **full copies** of reference tables

Scale vertically and split based on load

Perform **local planning** and **execution** of query fragments

Execute local transaction logic

Backed by Aurora distributed storage

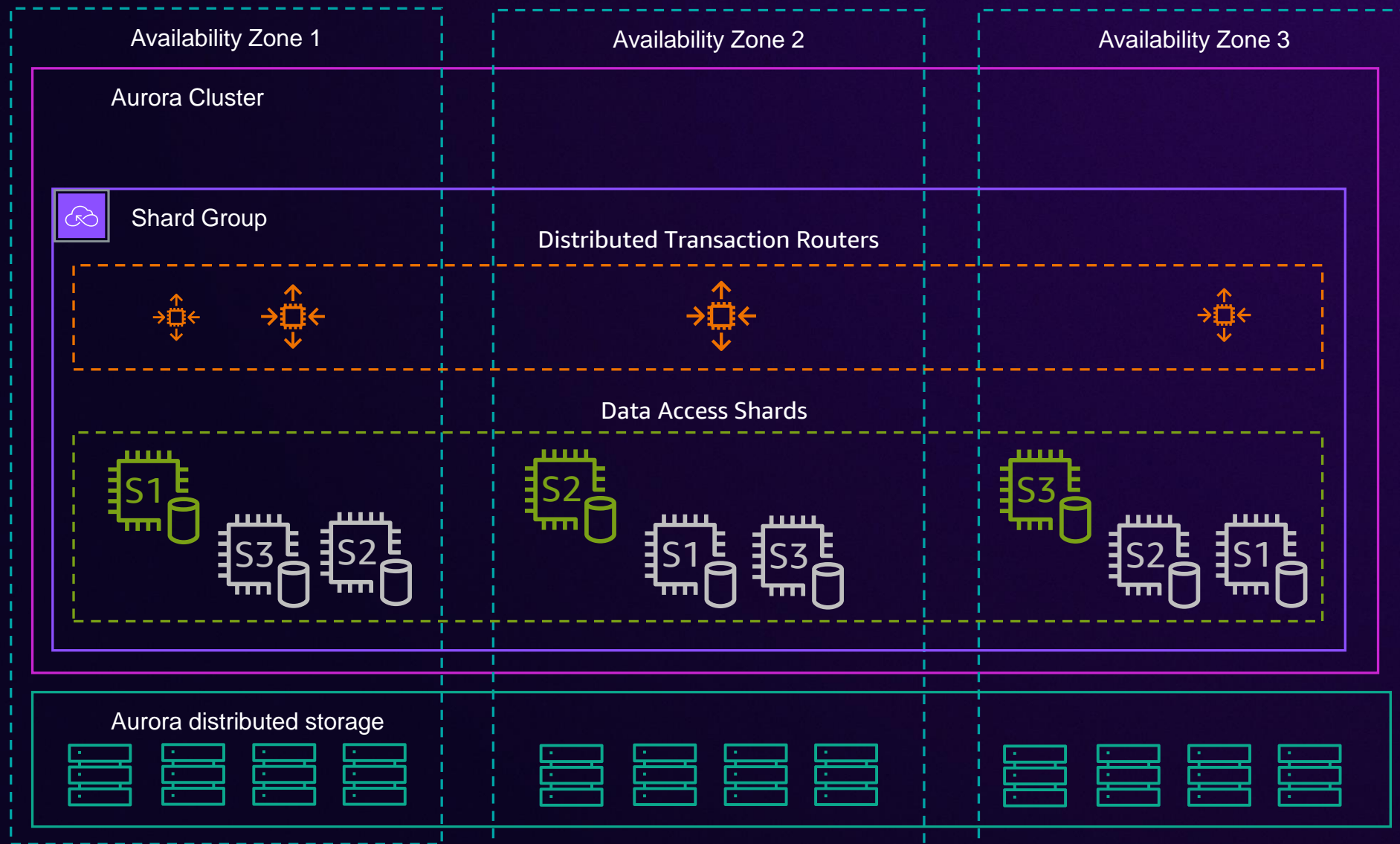


Topology and availability

Topology distributed across Availability Zones

Compute redundancy 0, 1, and 2

Routers are fungible and can be replaced as necessary



Data distribution



Sharded tables

```
SET rds_aurora.limitless_create_table_mode='sharded';
SET rds_aurora.limitless_create_table_shard_key='{cust_id}';
```

```
CREATE TABLE customer (
    cust_id INT PRIMARY KEY NOT NULL,
    name TEXT,
    email VARCHAR(100)
);
```

```
postgres_limitless=> \d+ customer
```

Partitioned table "public.customer"

Column	Type	Collation	Nullable	Default	Storage	Compression	Stats target
Description							
-----+-----+-----+-----+-----+-----+-----+-----							

cust_id	integer		not null		plain		
name	text				extended		
email	char..var(100)				extended		

```
Partition key: HASH (cust_id)
Partitions: customer_fs00001 FOR VALUES FROM (MINVALUE) TO ('-4611686018427387904'),
            customer_fs00002 FOR VALUES FROM ('-4611686018427387904') TO ('0'),
            customer_fs00003 FOR VALUES FROM ('0') TO ('4611686018427387904'),
            customer_fs00004 FOR VALUES FROM ('4611686018427387904') TO (MAXVALUE)
```

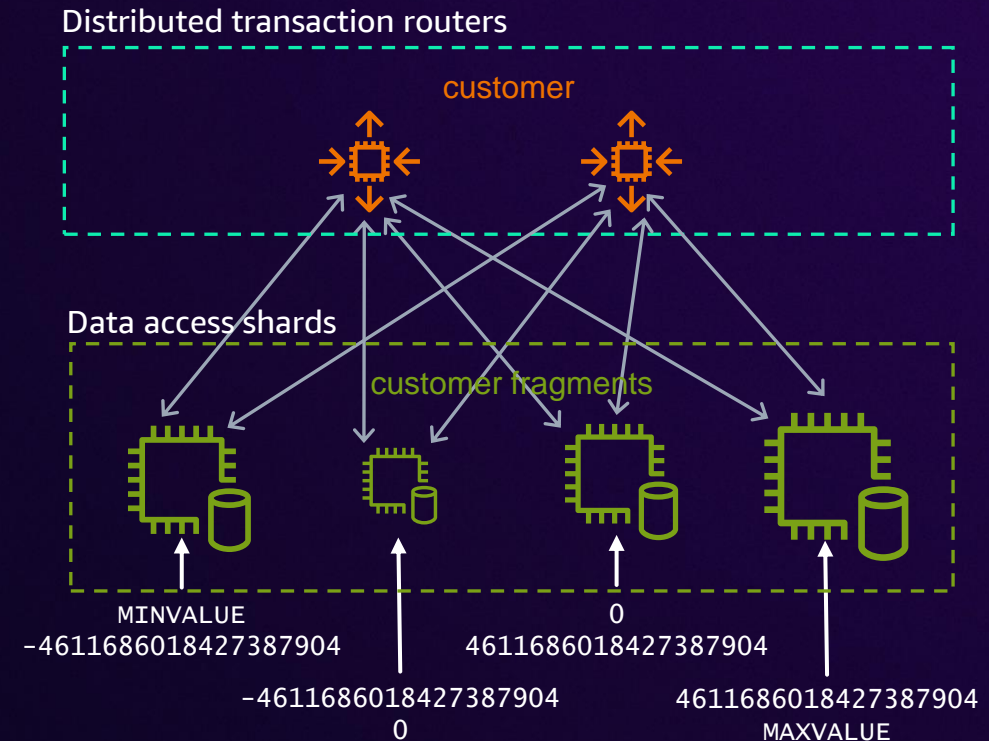

Hash-range partitioning

Shard key is **hashed** to 64-bits

Ranges of **64-bit** space are assigned to shards

Shards own **table fragments**

Routers have table fragment references, but **no data**



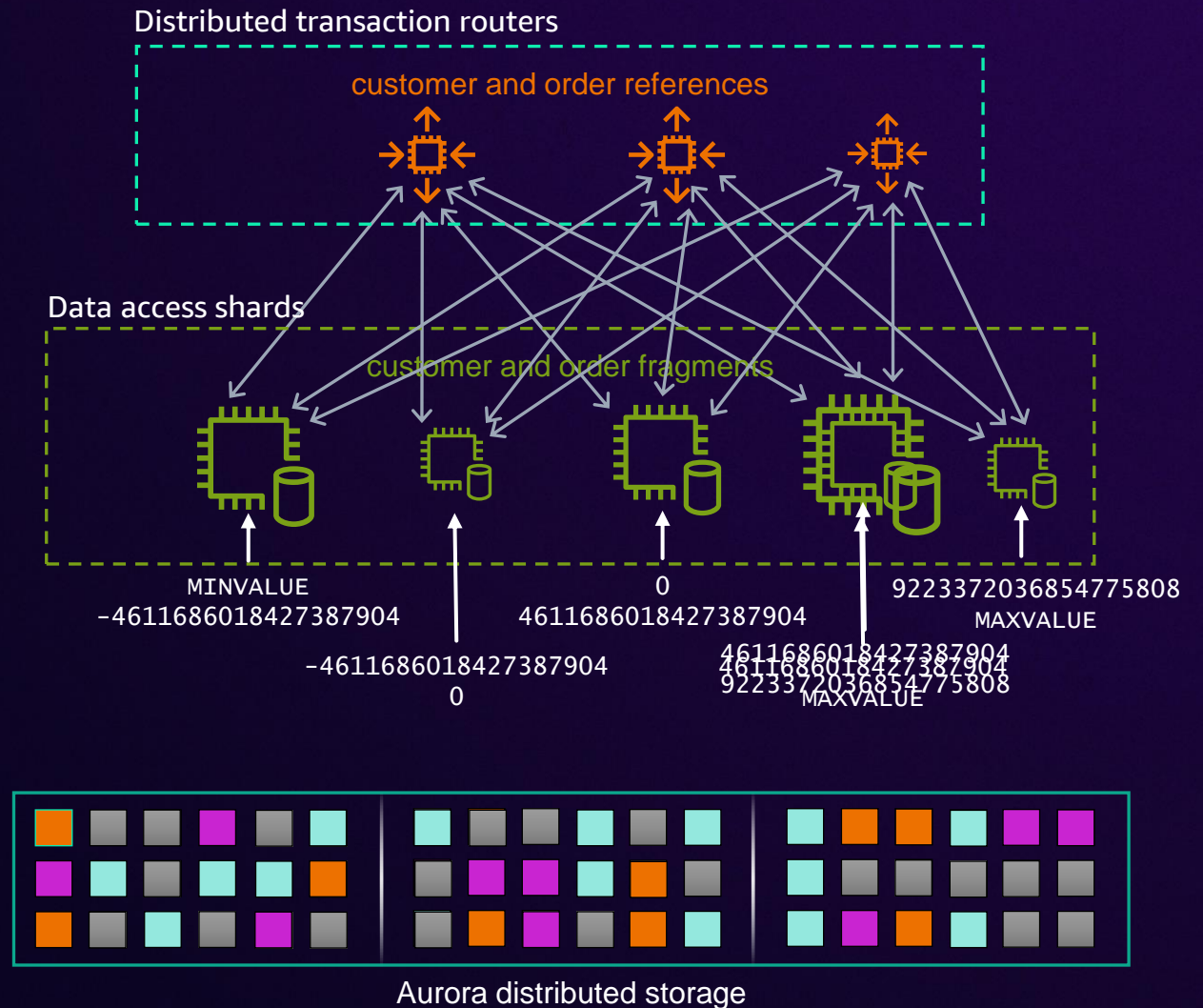
Horizontal scale out

“Shard split” occurs due to utilization or storage size

Collocated key ranges are moved together

Leverages Aurora storage level cloning and replication

Routers can be added



Reference tables

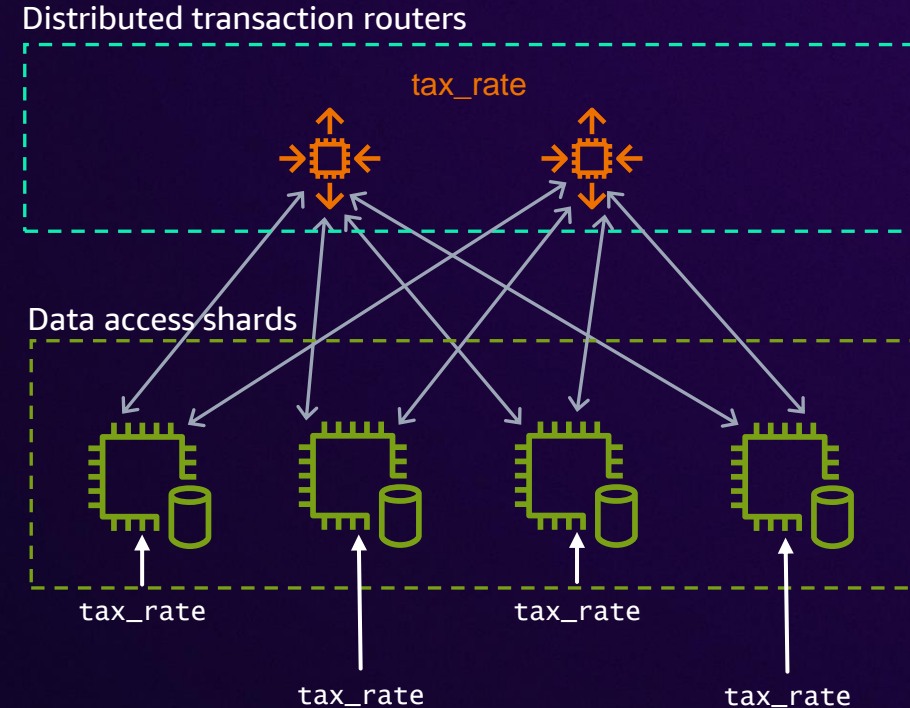
```
SET rds_aurora.limitless_create_table_mode='reference';
```

```
CREATE TABLE tax_rate (  
  tax_rate_id INT PRIMARY KEY NOT NULL,  
  city        TEXT NOT NULL,  
  state       TEXT,  
  country     TEXT NOT NULL,  
  tax_rate    DOUBLE NOT NULL  
);
```

Strongly consistent (ACID writes)

Enables join pushdown

Frequent read/join, infrequent write



Transactions



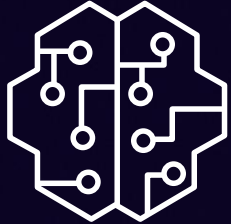
© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.

Transaction design goal

PostgreSQL semantics for **READ COMMITTED** and **REPEATABLE READ**

...with a consistent view as in a single system

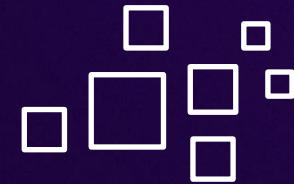
Challenges in a distributed database



Coordination limits
scalability



Transaction scope
unknown until commit



Query fragments execute
at different times



Maintain order



Consistent restores

Limitless ACID properties

PostgreSQL *read committed, repeatable read (SI)*

- Same semantics as single node PostgreSQL
- External consistency

Multi-shard writes are atomic

- All participants commit at same effective time

DDLs are transactional and strongly consistent

- DDLs are always RC, same as PostgreSQL
- Can be initiated from any SQL connection

Single system backup

- Point-in-time-restore fully consistent

PostgreSQL MVCC snapshots simplified

PostgreSQL snapshots are taken as of *now*

- Knows the current XID . . . a transaction start sequence number
- Records the XID of all running transactions
- Compares this with XID recorded on the tuples it is scanning
 - Tuple header has XID of inserter, deleter

Illustrative only - full implementation is more complex

PostgreSQL snapshots simplified

Invisible Tuples

- Written by an XID > snapshot XID
- Written by an XID that isn't committed
- Written by an XID running when the snapshot was taken

Visible Tuples

- Your own uncommitted writes
- Everything else

Illustrative only - full implementation is more complex

Limitless does snapshots as of *then*

Transaction router establishes the snapshot time *now*

Router passes this time to shards along with query fragment

Shards create their local snapshot as of *then*

Multi-shard snapshots will use the same time on all shards...*then*

Is transaction commit time earlier than snapshot time? Visible.

Distributed clocks

Historically, relying on wall time between multiple systems doesn't work

Innovation and major investment in time infrastructure makes this possible in AWS

Database algorithms built on highly reliable, drift bounded clocks

Extremely scalable design

Bounded clocks in EC2

Amazon Time Sync Service provides high quality time to EC2 instances

ClockBound is an open source daemon that provides *{earliest,latest}* uncertainty bounds, typically $< 1\text{msec}$

Actual true time guaranteed between *{earliest,latest}*

New architecture has clock source on Nitro card, $< 50\text{ usec}$ uncertainty

Repeatable read – distributed (with clocks)

Transaction T1

- 1) router gets time **t100**
- 2) execute on shard w/cust_id **619** using snapshot@**t100**

```
BEGIN TRANSACTION ISOLATION LEVEL  
REPEATABLE READ;
```

```
SELECT status FROM order WHERE cust_id  
= 619 and order_id = 61890340;  
filling
```

- 1) execute on shard w/cust_id **801** using snapshot@**t100**

```
SELECT status FROM order WHERE cust_id  
WHERE cust_id = 801 and order_id =  
80044011;  
filling
```

Transaction T2

- 1) router gets time **t103**
- 2) execute on shard w/cust_id **801** using snapshot@**t103**

```
BEGIN;
```

```
SELECT status FROM order WHERE cust_id  
WHERE cust_id = 801 and order_id =  
80044011;
```

filling

```
UPDATE order  
cust_id = 801;  
COMMIT;
```

- 1) router uses 1PC on shard
- 2) shard assigns commit@**t110**
- 3) acks commit when
 - a) writes durable on disk
 - b) earliest possible time > **t110**

Transaction T3

```
SELECT status FROM order WHERE cust_id  
WHERE cust_id = 801 and order_id =  
80044011;
```

shipped

- 1) router gets time **t125**
- 2) execute on shard w/cust_id **801** using snapshot@**t125**

Multi-shard writes

Build on modified two-phase commit protocol

Router coordinates distributed commit

All shards will commit the transaction with the same commit time

Commit latency is roughly 2–3x single shard commit

Transactions conclusion

Same RC/RR semantics as PostgreSQL

All reads are consistent, w/o quorum, even on failover

Commits w/single shard writes scale linearly (millions/sec)

Distributed commits are atomic

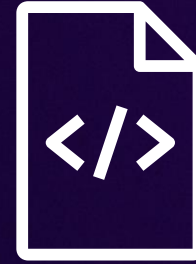
Queries & Performance



Fundamentally Aurora PostgreSQL



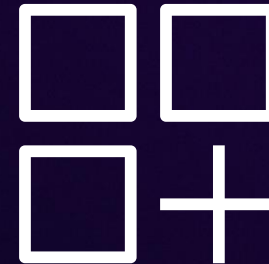
PostgreSQL wire compatible



PostgreSQL parser and semantics

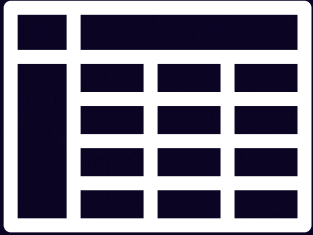


Broad surface area coverage



Selected extensions

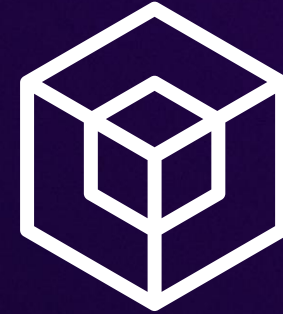
Query execution basics



PostgreSQL foreign tables
foundation



Enhancements in core engine



A custom foreign data wrapper

Query flow

Router

1. Receives query from client
2. Plans what can be sent to shards and any joins that must be done
3. Sends partial queries to shards with transaction context

7. Router does final joins, filters, and aggregations as necessary

Shard

4. Receives partial query from router
5. Plans local joins and scans
6. Execute and sent results to router

Locality is key to performance

Lowest latency and best scalability when locality is maintained

Push execution close to the data

Reduce messaging, leverages caching

Collocated and reference tables are key building blocks

Single shard optimization

Best performance when router determines query goes to a single shard

```
postgres_limitless=> EXPLAIN (VERBOSE, COSTS OFF) SELECT * FROM customers WHERE  
customer_id = 100;
```

QUERY PLAN

Foreign Scan

Output: customer_id, other_id, customer_name, balance

Remote SQL: SELECT customer_id,

other_id,

customer_name,

balance

FROM public.customers

WHERE (customer_id = 100)

Single Shard Optimized

Single shard join pushdown

```
postgres_limitless=> EXPLAIN (VERBOSE, COSTS OFF) SELECT * FROM orders  
LEFT JOIN zipcodes ON orders.zipcode_id = zipcodes.zipcode_id  
WHERE customer_id = 11;
```

QUERY PLAN

Foreign Scan

Output: customer_id, order_id, zipcode_id, customer_name, balance, zipcodes.zipcode_id, zipcodes.city

Remote SQL: SELECT orders.customer_id,
orders.order_id,
orders.zipcode_id,
orders.customer_name,
orders.balance,
zipcodes.zipcode_id,
zipcodes.city

FROM (public.orders
LEFT JOIN public.zipcodes ON ((orders.zipcode_id = zipcodes.zipcode_id)))
WHERE (orders.customer_id = 11)

single shard optimized

Function distribution

Collection of statements that operate on the same key value can be wrapped in a function

Significant improvement in latency and reduction in router CPU

See “Function distribution” and `Limitless_distribute_function` in the docs

Parallel operations

Parallel operations speed up via multi-shard execution

Some examples:

- Create index

- Analyze

- Vacuum

- Aggregates (count, sum, min, max)

Low latency at scale

Experiment in us-east-1

Three big client drivers
distributing random
updates across 100B rows

```
[screen 1: ec2-user@ip-172-31-36-202:/wip/pgbench-kermi15]
pgbench (16.4)
progress: 60.0 s, 686401.9 tps, lat 2.357 ms stddev 0.639, 0 failed
progress: 120.0 s, 849197.9 tps, lat 2.355 ms stddev 0.639, 0 failed
progress: 180.0 s, 848990.5 tps, lat 2.356 ms stddev 0.660, 0 failed
progress: 240.0 s, 849730.9 tps, lat 2.354 ms stddev 0.649, 0 failed
progress: 300.0 s, 836900.8 tps, lat 2.390 ms stddev 2.488, 0 failed
progress: 360.0 s, 847194.7 tps, lat 2.361 ms stddev 1.471, 0 failed
progress: 420.0 s, 848091.6 tps, lat 2.358 ms stddev 0.651, 0 failed
progress: 480.0 s, 846166.1 tps, lat 2.361 ms stddev 0.651, 0 failed
progress: 540.0 s, 850282.2 tps, lat 2.358 ms stddev 0.651, 0 failed
progress: Infinity s, 0.0 tps, lat 0.000 ms stddev 0.000, 0 failed
transaction type: ./limitless_pgbench_update.sql
scaling factor: 1000000
query mode: simple
number of clients: 2000
number of threads: 2000
maximum number of tries: 1
duration: 600 s
number of transactions actually processed: 848475
number of failed transactions: 0 (0.000%)
latency average = 2.355 ms
latency stddev = 1.228 ms
initial connection time = 11526.511 ms
tps = 848475.083080 (without initial connection time)
$:
```

```
[screen 1: ec2-user@ip-172-31-7-182:/wip/pgbench-kermi15]
pgbench (16.4)
progress: 60.0 s, 732174.1 tps, lat 2.206 ms stddev 0.701, 0 failed
progress: 120.0 s, 897809.3 tps, lat 2.228 ms stddev 0.675, 0 failed
progress: 180.0 s, 896915.5 tps, lat 2.230 ms stddev 0.686, 0 failed
progress: 240.0 s, 898754.1 tps, lat 2.225 ms stddev 0.674, 0 failed
progress: 300.0 s, 883947.6 tps, lat 2.262 ms stddev 2.486, 0 failed
progress: 360.0 s, 893172.1 tps, lat 2.239 ms stddev 1.469, 0 failed
progress: 420.0 s, 895048.2 tps, lat 2.234 ms stddev 0.674, 0 failed
progress: 480.0 s, 892966.3 tps, lat 2.240 ms stddev 1.598, 0 failed
progress: 540.0 s, 897730.2 tps, lat 2.228 ms stddev 0.681, 0 failed
progress: Infinity s, 0.0 tps, lat 0.000 ms stddev 0.000, 0 failed
transaction type: ./limitless_pgbench_update.sql
scaling factor: 1000000
query mode: simple
number of clients: 2000
number of threads: 2000
maximum number of tries: 1
duration: 600 s
number of transactions actually processed: 527475230
number of failed transactions: 0 (0.000%)
latency average = 2.231 ms
latency stddev = 1.211 ms
initial connection time = 11600.855 ms
tps = 896127.475507 (without initial connection time)
$:
```

```
[screen 1: ec2-user@ip-172-31-39-172:/wip/pgbench-kermi15]
-n -P 60 -T 600
pgbench (16.4)
progress: 60.0 s, 585688.0 tps, lat 2.451 ms stddev 0.611, 0 failed
progress: 120.0 s, 811500.4 tps, lat 2.464 ms stddev 0.612, 0 failed
progress: 180.0 s, 811900.0 tps, lat 2.463 ms stddev 0.627, 0 failed
progress: 240.0 s, 813461.0 tps, lat 2.458 ms stddev 0.619, 0 failed
progress: 300.0 s, 801854.7 tps, lat 2.494 ms stddev 2.531, 0 failed
progress: 360.0 s, 810100.2 tps, lat 2.469 ms stddev 1.505, 0 failed
progress: 420.0 s, 812459.2 tps, lat 2.461 ms stddev 0.621, 0 failed
progress: 480.0 s, 809339.0 tps, lat 2.471 ms stddev 1.564, 0 failed
progress: 540.0 s, 813195.7 tps, lat 2.459 ms stddev 0.635, 0 failed
transaction type: ./limitless_pgbench_update.sql
scaling factor: 1000000
query mode: simple
number of clients: 2000
number of threads: 2000
maximum number of tries: 1
duration: 600 s
number of transactions actually processed: 472903957
number of failed transactions: 0 (0.000%)
latency average = 2.465 ms
latency stddev = 1.185 ms
initial connection time = 17012.864 ms
tps = 811034.386711 (without initial connection time)
$:
```

Get started today!

Aurora Limitless Database - *new* [Info](#)

With Limitless Database, Aurora can automatically scale write throughput and data storage capacity beyond the limits of a single DB cluster.

DB shard group identifier

Type a name for your DB shard group. The name must be unique across all DB shard groups owned by your AWS account in the current AWS Region.

Enter DB shard group identifier

Constraints: 1 to 60 alphanumeric characters or hyphens. First character must be a letter. Can't contain two consecutive hyphens. Can't end with a hyphen.

DB shard group capacity range [Info](#)

Enter the minimum and maximum capacity for Limitless Database. The capacity is measured in Aurora capacity units (ACUs) across all routers and shards.

Minimum capacity (ACUs)

24

(48 GiB)

Enter a value greater than or equal to 16 ACUs.

Maximum capacity (ACUs)

384

(768 GiB)

Enter a value less than or equal to 6144 ACUs.

DB shard group deployment

The number of additional cross Availability Zone standby shards. Adding compute redundancy will have a significant impact on cost. [Learn more](#)

- ☒ **No compute redundancy**
Creates a DB shard group without standbys for each shard.
- ☐ **Compute redundancy with a single failover target**
Each shard is created with one compute standby in a different Availability Zone.
- ☐ **Compute redundancy with two failover targets**
Each shard is created with two compute standbys in different Availability Zones.

aws rds **create-db-shard-group**

--db-cluster-identifier proddb

--db-shard-group-identifier proddb-sg

--min-acu 150

--max-acu 600

--compute-redundancy 2



Summary

- ✓ Challenges for scaling
- ✓ Scales to millions of write transactions per second
- ✓ Manages petabytes of data
- ✓ Scalable architecture
- ✓ Data distribution
- ✓ Query and transactions



**CALL TO
ACTION**

Get started today with AWS console: <https://console.aws.amazon.com/>

Learn more: <https://docs.aws.amazon.com/AmazonRDS/latest/AuroraUserGuide/limitless.html>

Related sessions

Session ID	Session Title
DAT416	Scalable database solutions with Aurora PostgreSQL Limitless Database
DAT316	Build scalable and cost-optimized apps with Amazon Aurora Serverless
DAT424	Get started with the latest Amazon Aurora innovations
DAT405	Deep dive into Amazon Aurora and its innovations
DAT304	Amazon Aurora HA and DR design patterns for global resilience



Thank you!

Anum Jang Sher
anujangs@amazon.com

David Wein
dcw@amazon.com



Please complete the session
survey in the mobile app