# Deciphering 2phase commit

Ashwin Agrawal
Asim Praveen
{aagrawal, apraveen}@pivotal.io

# Scale out

- Single instance is limited
- Manual attempts at sharding PostgreSQL
- FDW based sharding
- MPP → distributed databases

# Challenge with atomicity

```
begin;

insert into account values (id = 1 ...);

insert into account values (id = 2 ...);

commit;
```

shard 1

shard 2

# Challenge with atomicity
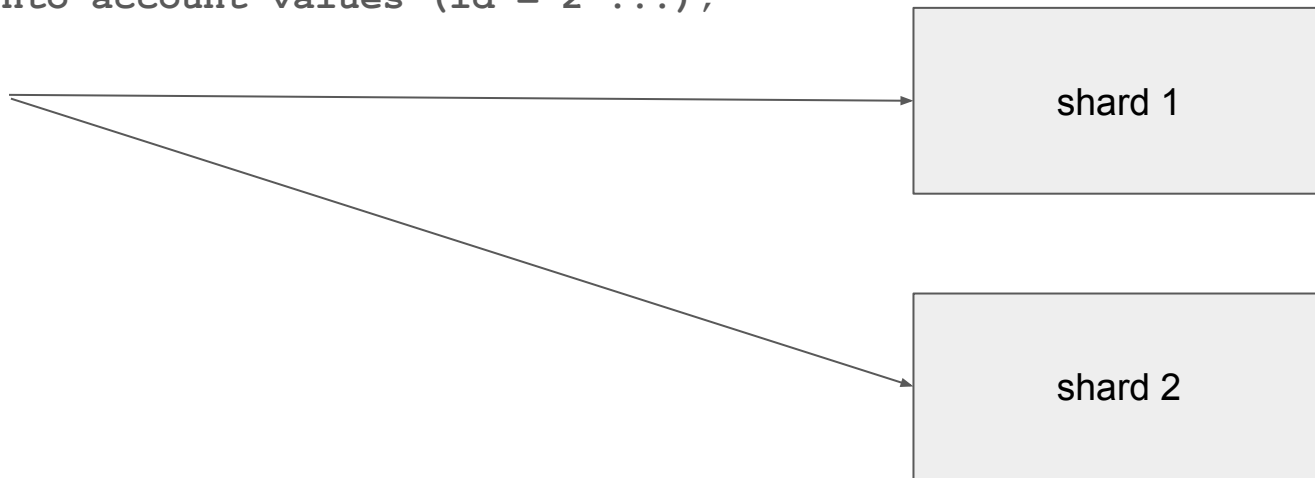
```
begin;

insert into account values (id = 1 ...);

insert into account values (id = 2 ...);

commit;
```

shard 1

shard 2

# Challenge with atomicity

```
begin;

insert into account values (id = 1 ...);

insert into account values (id = 2 ...);

commit;
```

shard 1

shard 2

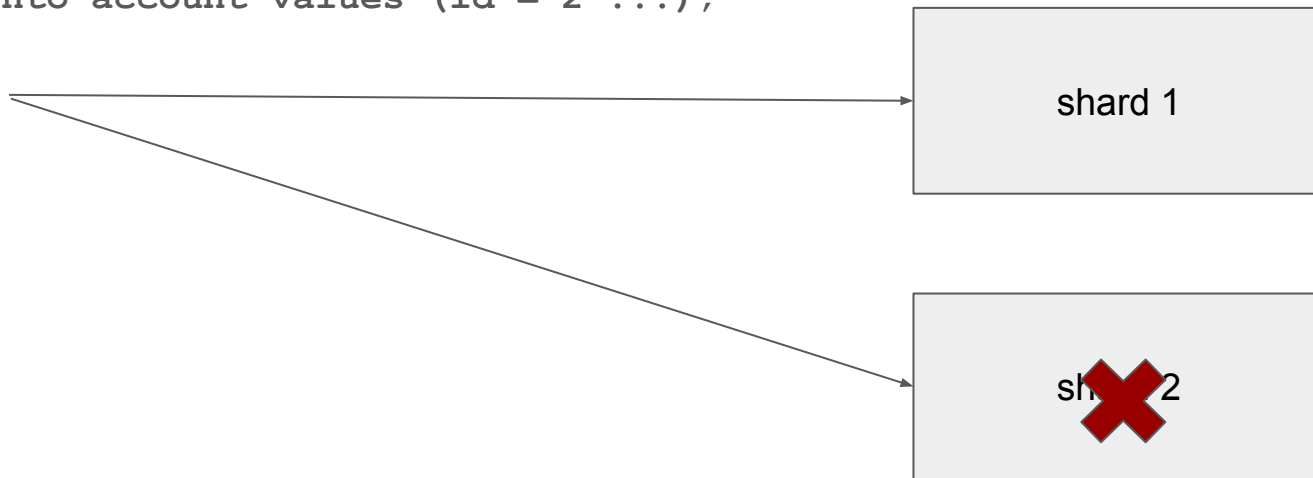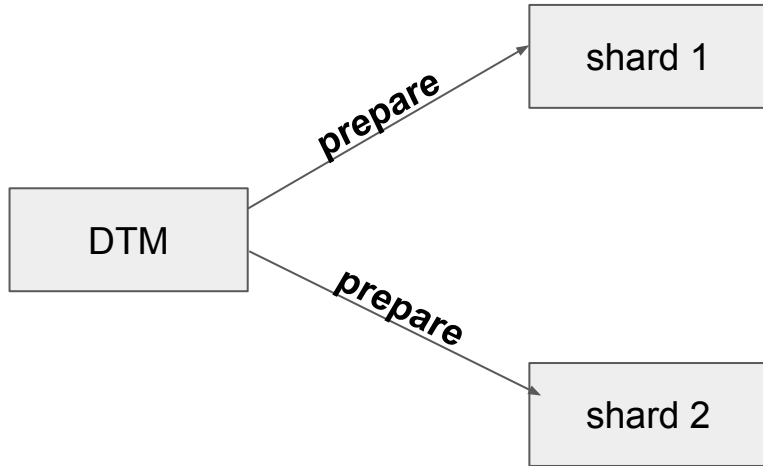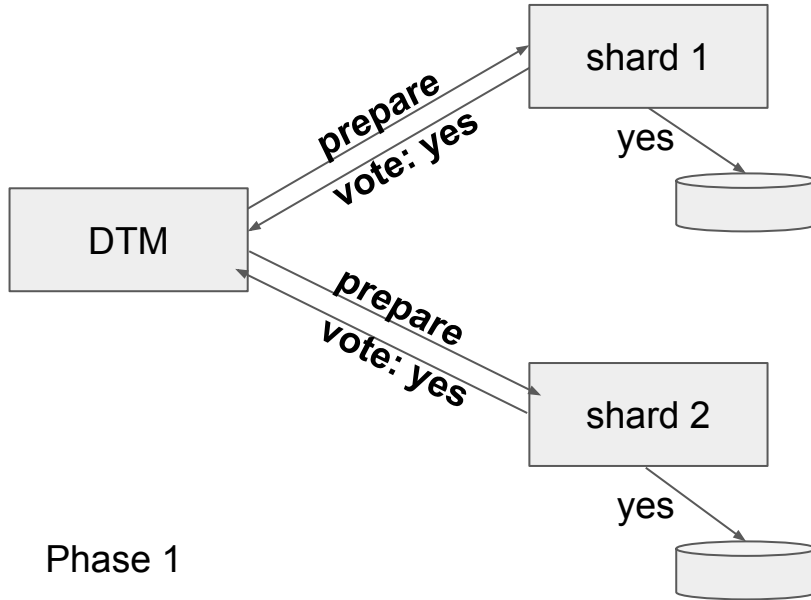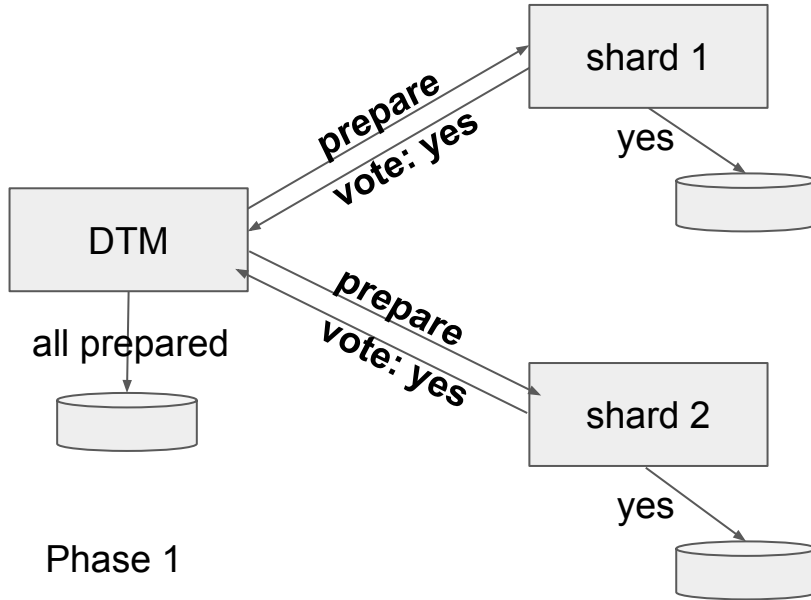# Two phase commit



Phase 1

# Two phase commit



Phase 1

# Two phase commit



shard 1

prepare
vote: yes

yes

DTM

all prepared

prepare
vote: yes

shard 2

yes

Phase 1

# Two phase commit



Phase 1

Phase 2

# 2PC: shard failure in phase one



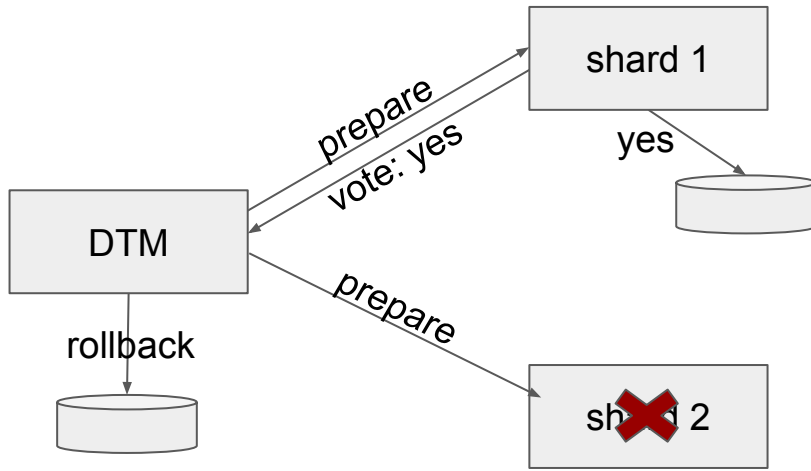prepare

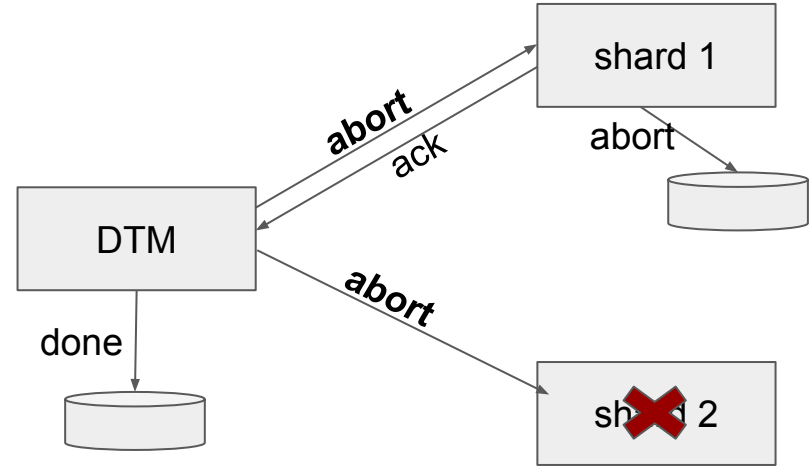vote: yes

shard 1

yes

DTM

rollback

prepare

shard 2

Phase 1

# 2PC: shard failure in phase one



Phase 1

Phase 2

# 2PC: shard failure in phase 2



prepare
vote: yes

yes

all prepared

prepare
vote: yes

yes

Phase 1

commit
ack

commit

commit

Phase 2

# 2PC: recovery of shard 2

timeline

DTM

Ack not received from shard 2;

retry sending commit to shard 2

shard 2

checkpoint

prepared: xid

wait for commit/abort message

# 2PC: DTM crashed in phase 1

shard 1

prepare

vote

DTM ✖

prepare

shard 2

vote

Phase 1

# 2PC: DTM crashed in phase 1



Phase 1

DTM Recovery

# 2PC: DTM crashed after phase 1



prepare
vote: yes

shard 1

yes

DTM ✖

all prepared

prepare
vote: yes

shard 2

yes

Phase 1

# 2PC: DTM crashed after phase 1



Phase 1

DTM Recovery

# 2PC vs 1PC

- Prepare phase
  - 1 network round trip
  - 1 disk flush
- Commit phase
  - 1 network round trip
  - 1 disk flush
- 2PC guarantees A and D of ACID

# Single node snapshot isolation

Tuple headers contain:

- xmin: transaction ID of inserting transaction

- xmax: transaction ID of replacing/deleting transaction (initially NULL)

Basic idea: tuple is visible if xmin is **valid** and xmax is not. "**Valid**" means "**either committed or the current transaction**".

# "Snapshot" filter away active transactions

Rules ensuring no transaction committing after the current transaction's start be considered committed:

- Currently running transactions IDs never considered valid, even if shown committed in pg_clog.
- Transaction ID higher than the current transaction is not valid (future transaction).

# Challenge with isolation

```
A: begin;
B: begin;

A: insert into acc values(id=1, ...);
B: insert into acc values(id=3, ...);
```

shard1
  A: 10
  B: 15

shard2

# Challenge with isolation

```
A: begin;
B: begin;

A: insert into acc values(id=1, ...);
B: insert into acc values(id=3, ...);

B: insert into acc values(id=4, ...);
A: insert into acc values(id=2, ...);

B: commit;
```

shard1
A: 10
B: 15

shard2
B: 20
A: 25

# Challenge with isolation

```
A: begin;
B: begin;

A: insert into acc values(id=1, ...);
B: insert into acc values(id=3, ...);

B: insert into acc values(id=4, ...);
A: insert into acc values(id=2, ...);

B: commit;
A: select * from acc;
    1, 2, 4
```

shard1
    A: 10
    B: 15

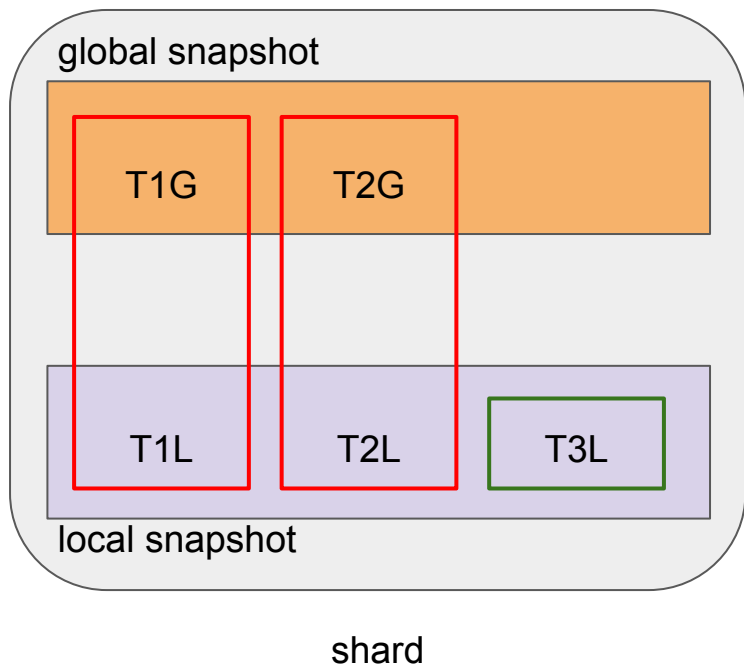B is in future for A

shard2
    B: 20
    A: 25

B visible to A !!!

# Global snapshot



shard

- Global xid and global snapshot provided by DTM
- Gxmin, Gxmax, gInProgress [ ]
- tuples contain local xmin/xmax

```
if (!XidInSnapshot(GS, xid))
{
    XidInSnapshot(LS, xid)
}
```
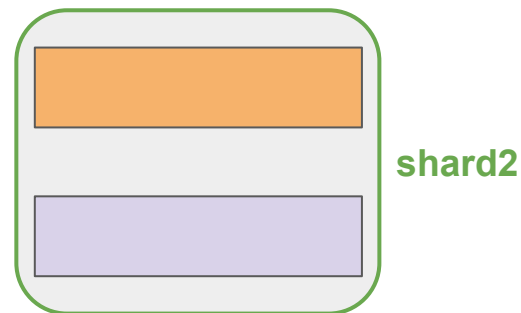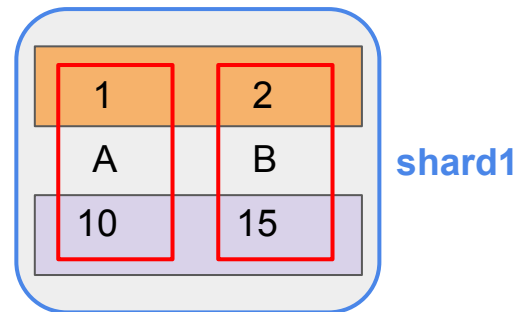
# Global snapshot in action

```
A: begin; GXID: 1
B: begin; GXID: 2

A: insert into acc values(id=1, ...);
B: insert into acc values(id=3, ...);
```

# Global snapshot in action

```
A: begin; GXID: 1
B: begin; GXID: 2

A: insert into acc values(id=1, ...);
B: insert into acc values(id=3, ...);

B: insert into acc values(id=4, ...);
A: insert into acc values(id=2, ...);

B: commit;
```

# Global snapshot in action
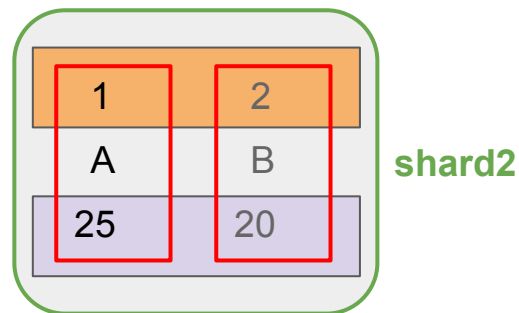
```
A: begin; GXID: 1
B: begin; GXID: 2

A: insert into acc values(id=1, ...);
B: insert into acc values(id=3, ...);

B: insert into acc values(id=4, ...);
A: insert into acc values(id=2, ...);

B: commit;
A: select * from acc;
    1, 2
```

# Global snapshot with local transaction

```
A: begin; GXID: 1
B: begin; GXID: 2

A: insert into acc values(id=1, ...);
B: insert into acc values(id=3, ...);

B: insert into acc values(id=4, ...);
A: insert into acc values(id=2, ...);
L: select * from acc;
       /* 0 rows */
B: commit;
L: select * from acc;
     4
```
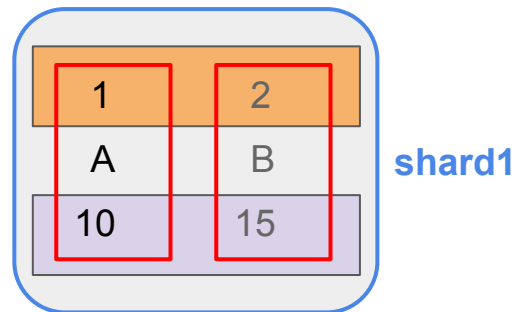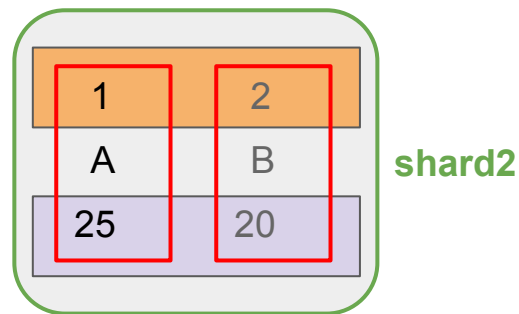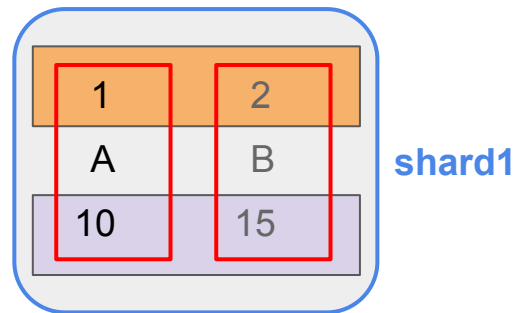
# Implementation options

- Model
  - Pull
    - DTM as a service
    - participants join a transaction
    - transaction can be initiated by any participant
  - Push
    - DTM initiates transaction and decides participants
    - MPP databases

# Implementation options

- Transactions
  - Global and local
    - global and local transaction IDs, snapshots
    - mapping between global and local transaction IDs
  - Global only
    - only one xid and snapshot across the cluster

# ACID distributed system

- **Two phase commit ⇒ A and D**
- **Global snapshot ⇒ I**

**?? Spot the Problem ??**

```
A: begin;
A: update acc set ... where id = 1
B: begin;
B: update acc set ... where id = 2
B: update acc set ... where id = 1
A: update acc set ... where id = 2
```

# ACID distributed system

- **Two phase commit ⇒ A and D**
- **Global snapshot ⇒ I**

**?? Spot the Problem ??**

```
A: begin;
A: update acc set ... where id = 1
B: begin;
B: update acc set ... where id = 2
B: update acc set ... where id = 1
A: update acc set ... where id = 2
```

- Global lock manager ⇒ C

рахмат

danke

謝謝

ngiyabonga

teşekkür ederim

Баярлалаа

спасибо

faafetai lava

mersi

kia ora

barka

welalin

tack

spas

mahalo

tapadh leat

vinaka

dank je

misaotra

matondo

paldies

grazzi

kiitos

dankie

nandri

спасибі

blagodaram

хвала

asante

manana

nani

dhanyavad

köszönöm

akun

dankon

ačiū

mauruuru

obrigada

tenki

enkosi

bedankt

bayarlalaa

grâce

hvala

djere dieuf

tau

mochchakkeram

murakoze

gràcie

thank you

gracias

chokrane

dziękuję

chnorakaloutioun

gratias ago

gràcies

sulpáy

go raibh maith agat

mamnun

sobodi

dékuji

sagolun

kop khun krap

ţaiku

arigatô

takk

dakujem

trugarez

obrigado

mesi

sukriya

ありがとう

дякую

najis tuke

grazie

tanemirt

мерси

didi madloba

kam sah hamnida

terima kasih

rahmet

diolch

dhanyavadagalu

shukriya

merce

rahmat

তোমাকে ধন্যবাদ

감사합니다

xiexie

ευχαριστώ

merci