

A person is silhouetted against a bright, hazy sunset sky, sitting on the edge of a dark, rocky cliff. They are holding a laptop. The sun is low on the horizon, creating a strong orange and yellow glow that fills the lower half of the image. The sky is filled with soft, white clouds. The overall mood is contemplative and serene.

# Текущее состояние DTM

Стас Кельвич,  
Константин Книжник,  
Константин Пан

# Задача

- ▶ ACID-кластер

## Пример

- ▶ таблица accounts (id, amount)
- ▶ разбита по id на три сервера

Переводим деньги с одного счёта на другой.  
Возможно, с одного сервера на другой.

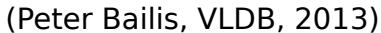
# Атомарность

COMMIT происходит либо на всех участвующих узлах, либо ни на одном.

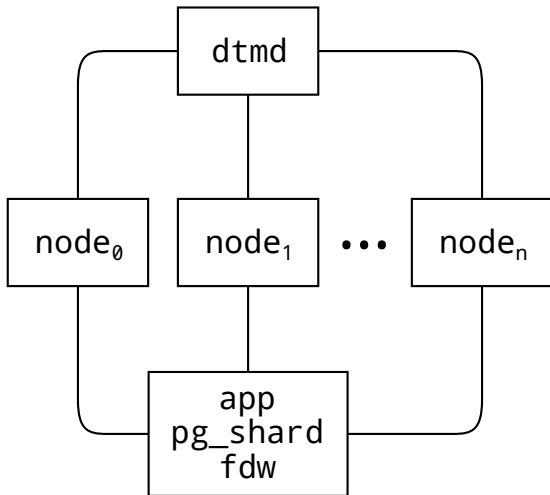
```
begin;  
update acc set a = a -  
100 where id = 1;  
select sum(a) from acc;  
  
update acc set a = a +  
100 where id = 2;  
commit;
```

# Уровни изоляции

Isolation Level	Dirty Read	Nonrepeatable Read	Phantom Read	Serialization Anomaly
Read uncommitted	Allowed, but not in PG	Possible	Possible	Possible
Read committed	Not possible	Possible	Possible	Possible
Repeatable read	Not possible	Not possible	Allowed, but not in PG	Possible
Serializable	Not possible	Not possible	Not possible	Not possible



# Структура кластера





## С точки зрения клиента

src server	dst server
<code>create extension pg_dtm;</code>	<code>create extension pg_dtm;</code>
<code>select begin_glob...();</code> <code>begin</code> <code>update ...;</code> <code>commit;</code>	<code>select join_glob...(xid);</code> <code>begin;</code>  <code>update ...;</code> <code>commit;</code>

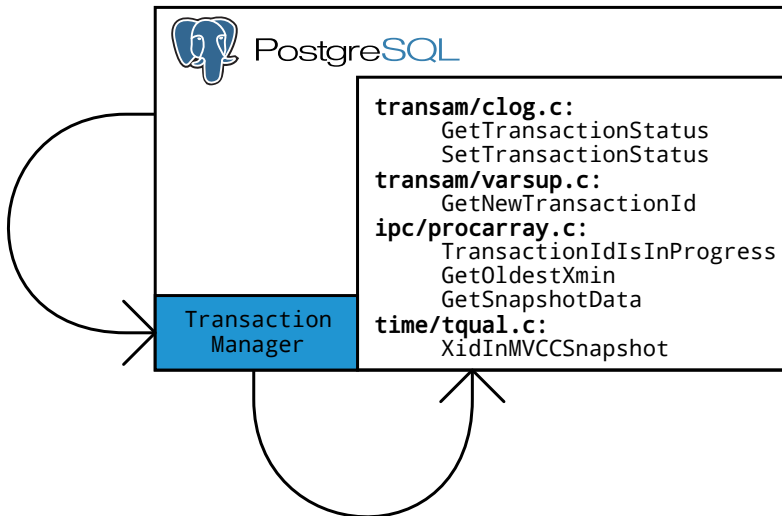
# С точки зрения кода: vanilla



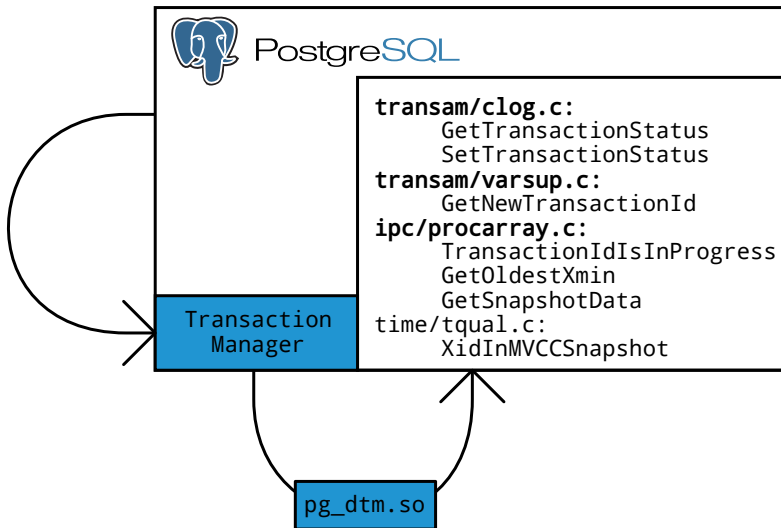
PostgreSQL

```
transam/clog.c:  
    GetTransactionStatus  
    SetTransactionStatus  
transam/varsup.c:  
    GetNewTransactionId  
ipc/proccarray.c:  
    TransactionIdIsInProgress  
    GetOldestXmin  
    GetSnapshotData  
time/tqual.c:  
    XidInMVCCSnapshot
```

## С точки зрения кода: после patch



# С точки зрения кода: после preload .so

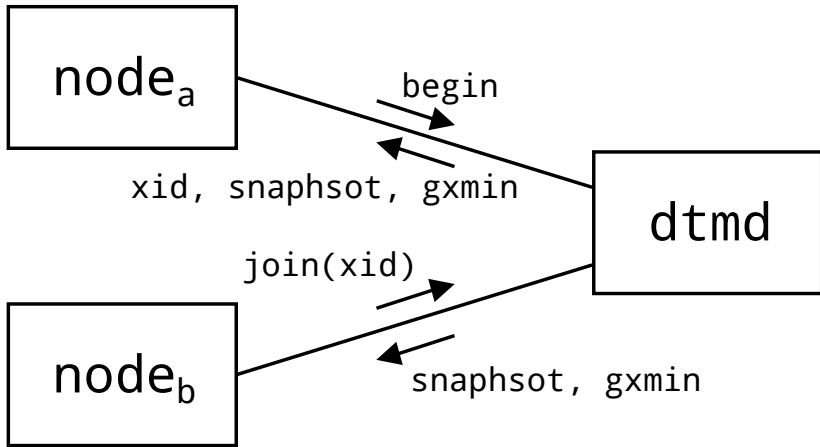


## С точки зрения кода: реализации транзакций

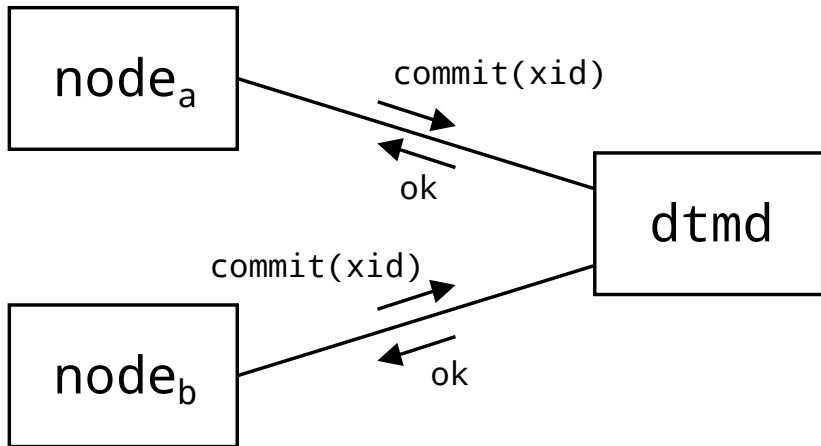
В качестве ТМ можно прикрутить разные реализации:

- ▶ **Snapshot Sharing** (XL, DTM)  
no 2pc, no local, daemon
- ▶ **Timestamp** (Spanner, Cockroach)  
2pc, local, no daemon
- ▶ **Incremental** (SAP HANA)  
no 2pc, locals, increments

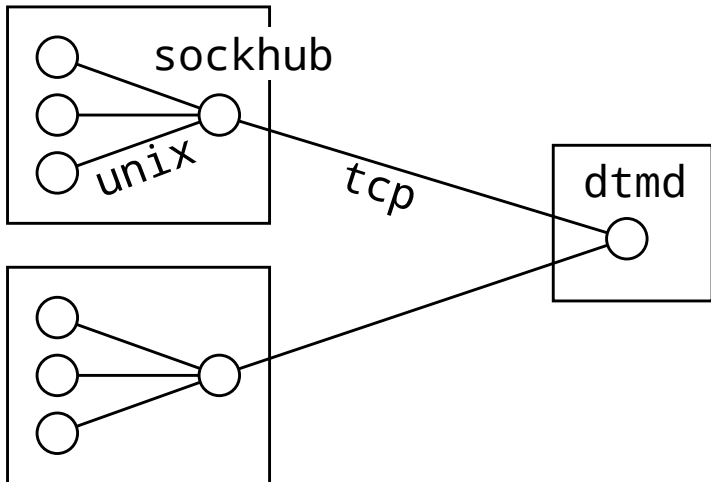
## Протокол (начало)



## Протокол (конец)



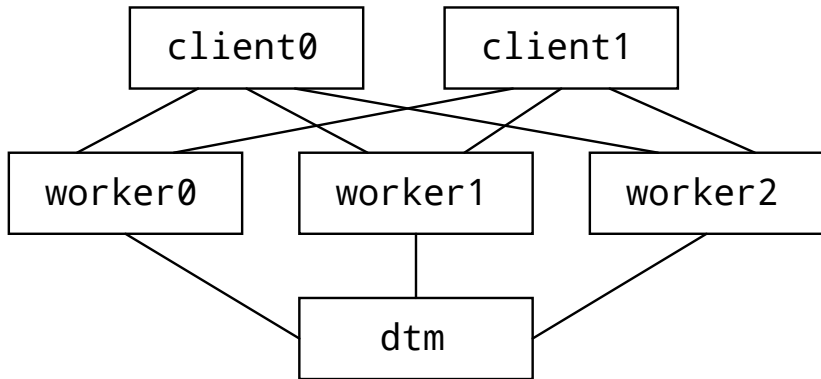
backends



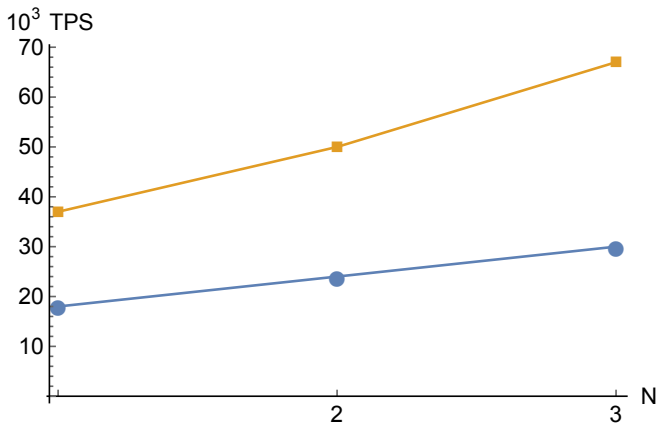


# Тесты

# Конфигурация



## Результаты



желтый — без DTM  
синий — с DTM

## Результаты

- ▶ pg\_shard: 17 k TPS
- ▶ fdw: 15 k TPS

# **Ограничения**

## Переключение контекста

На этом тесте без DTM упираемся в Context Switches. Включение DTM = x2 переключений контекста на транзакцию. Т.е. в 2 раза меньше tps с сервера.

Что можно сделать? Писать в очередь в общей памяти и будет быстрее.

DTM однопоточный.

Пусть нужно 3k ops на транзакцию (оптимистично).

$3 \text{ GHz} / 3 \text{ k.ops} = 1\text{M TPS}$

1 транзакция:  $112 + 8 \cdot N$  bytes

$N \approx 400$

около 1700 байт на транзакцию

- ▶ 1Gb: 300k tps
- ▶ 10Gb: 3M tps



# Репликация

...