

# Scaling out by distributing and replicating data in Postgres-XC

Ashutosh Bapat  
**@Postgres Open 2012**

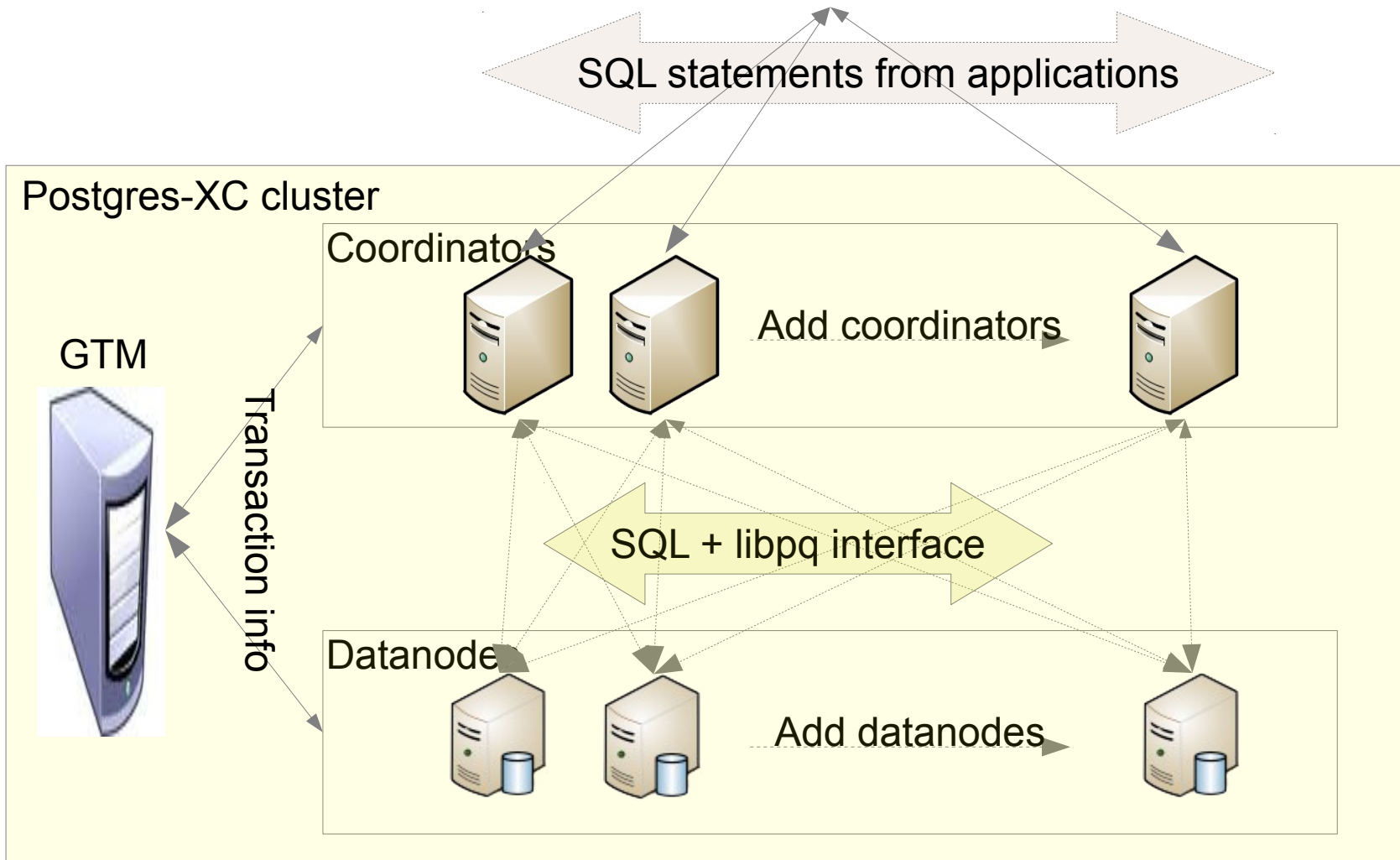
# Agenda

- What is Postgres-XC
- Postgres-XC architecture over-view
- Data distribution in XC
- Effect of data distribution on performance
- Example DBT-1 schema

# What is Postgres-XC

- **Shared Nothing Cluster**
  - Multiple collaborating PostgreSQL-like servers
  - No resources shared
  - Scaling by adding commodity hardware
- **Write-scalable**
  - Write/Read scalable by adding nodes
  - Multiple nodes where writes can be issued
- **Synchronous**
  - Writes to one node are reflected on all the nodes
- **Transparent**
  - Applications need not care about the data distribution

# Postgres-XC architecture



# Distribution strategies

- **Replicated tables**
  - Each row of the table is stored on all the datanodes where the table is replicated
- **Distributed tables**
  - Each row exists only on a single datanode
  - Distribution strategies
    - HASH
    - MODULO
    - ROUNDROBIN
    - User defined functions (TBD)

# Replicated Table

Writes



write write write



val	val2
1	2
2	10
3	4

val	val2
1	2
2	10
3	4

val	val2
1	2
2	10
3	4

Reads



read



val	val2
1	2
2	10
3	4

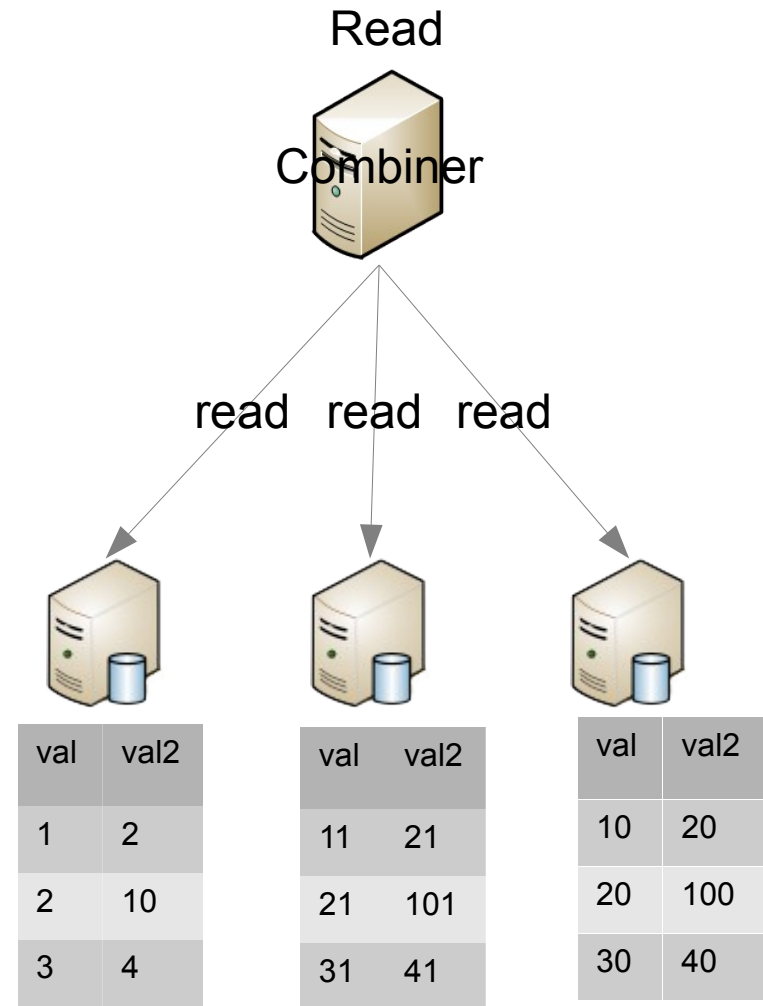
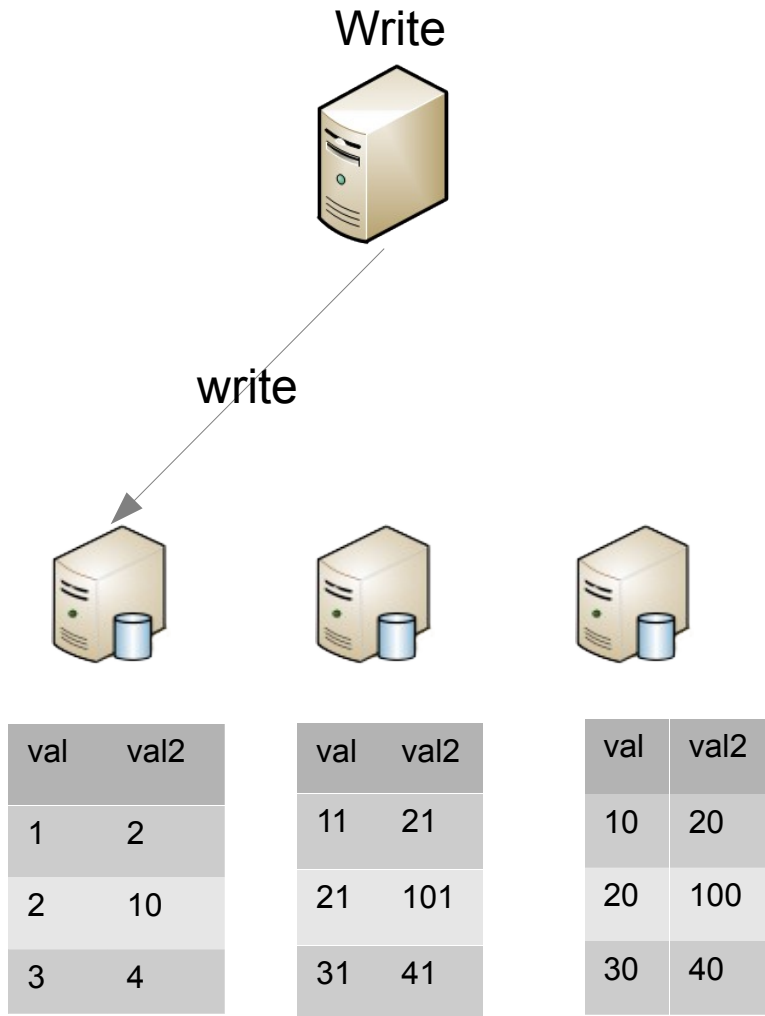
val	val2
1	2
2	10
3	4

val	val2
1	2
2	10
3	4

# Replicated Tables

- Statement level replication
- Each write needs to be replicated
  - writes are costly
- Read can happen on any node (where table is replicated)
  - reads from different coordinators can be routed to different nodes
- Useful for relatively static tables, with high read load

# Distributed Tables





# Distributed Tables

- Write to a single row is applied only on the node where the row resides
  - Multiple rows can be written in parallel
- Scanning rows spanning across the nodes (e.g. table scans) can hamper performance
- Point reads and writes based on the distribution column value show good performance
  - Datanode where the operation happens can be identified by the distribution column value



# Distributed query processing in Postgres- XC

# Distributed query processing in Postgres-XC

- **Coordinator**
  - Accepts queries and plans them
  - Finds the right data-nodes from where to fetch the data
  - Frames queries to be sent to these data-nodes
  - Gathers data from data-nodes
  - Processes it to get the desired result
- **Datanode**
  - Executes queries from coordinator like PostgreSQL
  - Has same capabilities as PostgreSQL

## Query processing balance

- Coordinator tries to delegate maximum query processing to data-nodes
  - Indexes are located on datanodes
  - Materialization of huge results is avoided in case of sorting, aggregation, grouping, JOINS etc.
  - Coordinator is freed to handle large number of connections
- Distributing data wisely helps coordinator to delegate maximum query processing and improve performance
- Delegation is often termed as shipping



# SQL prompt



# Deciding the right distribution strategy

## Read-write load on tables

- High point reads (based on distribution column)
  - Distributed or replicated
- High read activities but no frequent writes
  - Better be replicated
- High point writes
  - Better be distributed
- High insert-load, but no frequent update/delete/read
  - Better be round-robin

## Query analysis (Frequently occurring queries)

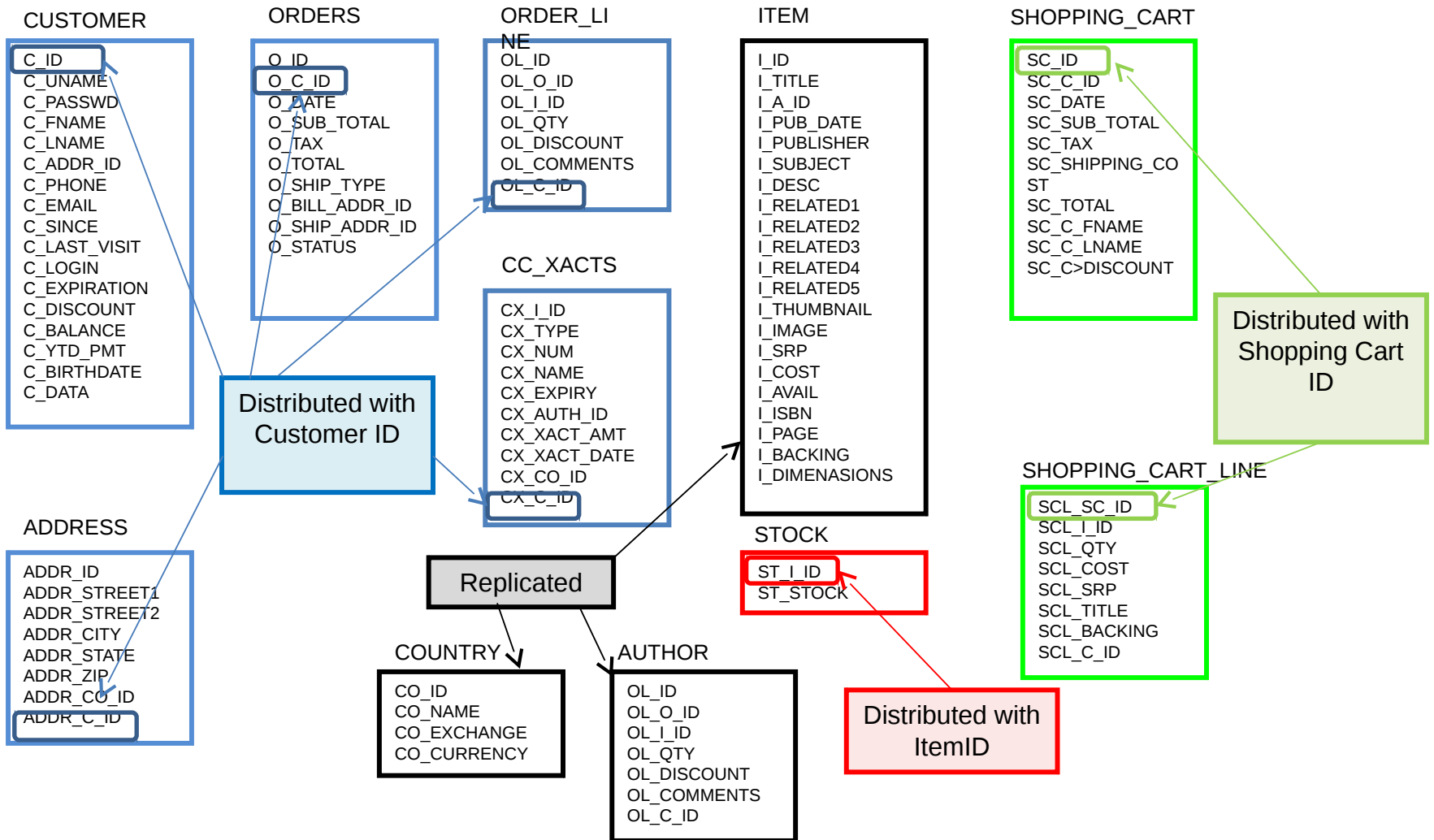
- Find the relations/columns participating in equi-Join conditions, WHERE clause etc.
  - Distribute on those columns
- Find columns participating in GROUP BY, DISTINCT clauses
  - Distribute on those columns
- Find columns/tables which are part of primary key and foreign key constraints
  - Global constraints are not yet supported in XC
  - Distribute on those columns



## Thumb rules

- Infrequently written tables participating in JOINS with many other tables (Dimension tables)
  - Replicated table
- Frequently written tables participating in JOINS with replicated tables
  - Distributed table
- Frequently written tables participating in JOINS with each other, with equi-JOINing columns of same data type
  - Distribute both of them by the columns participating in JOIN on same nodes
- Referenced tables
  - Better be replicated

# DBT-1 schema



## Example DBT-1 (1)

- author, item
  - Less frequently written
  - Frequently read from
  - Author and item are frequently JOINed
    - Dimension tables
  - Hence replicated on all nodes

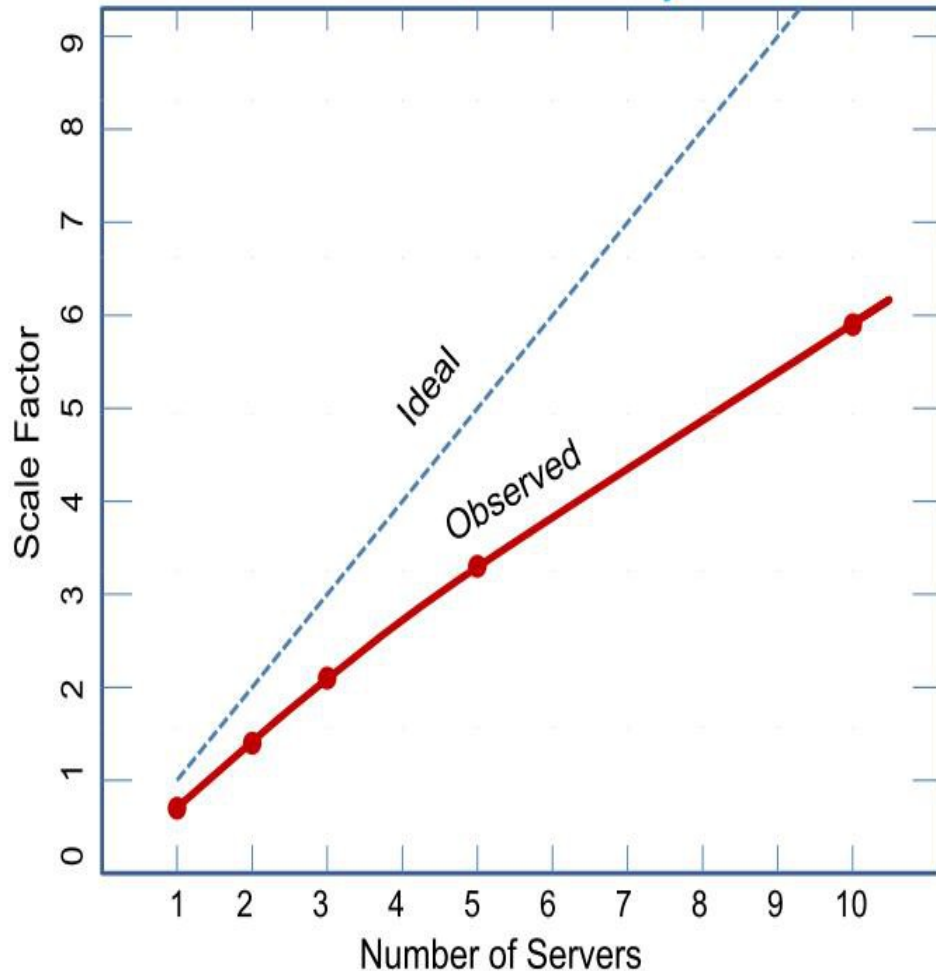
## Example DBT-1 (2)

- customer, address, orders, order\_line, cc\_xacts
  - Frequently written
    - hence distributed
  - Participate in JOINS amongst each other with customer\_id as JOIN key
  - point SELECTs based on customer\_id
    - hence distributed by hash on customer\_id so that JOINS are shippable
  - Participate in JOINS with item
    - Having item replicated helps pushing JOINS to datanode

## Example DBT-1 (3)

- **Shopping\_cart, shopping\_cart\_line**
  - Frequently written
    - Hence distributed
  - Point selects based on column shopping\_cart\_id
    - Hence distributed by hash on shopping\_cart\_id
  - JOINS with item
    - Having item replicated helps

# DBT-1 scale-up



- Old data, we will publish bench-marks for 1.0 soon.
- DBT-1 (TPC-W) benchmark with some minor modification to the schema
- 1 server = 1 coordinator + 1 datanode on same machine
- Coordinator is CPU bound
- Datanode is I/O bound



# Other scaling tips

## Using GTM proxy

- GTM can be a bottleneck
  - All nodes get snapshots, transactions ids etc. from GTM
- GTM-proxy helps reduce the load on GTM
  - Runs on each physical server
  - Caches information about snapshots, transaction ids etc.
  - Serves logical nodes on that server



# Adding coordinator and datanode

- **Coordinator**
  - Scaling connection load
  - Too much load on coordinator
  - Query processing mostly happens on coordinator
- **Datanode**
  - Data scalability
    - Number of tables grow – new nodes for new tables/databases
    - Distributed table sizes grow – new nodes providing space for additional data
  - Redundancy

# Impact of transaction management on performance

- 2PC is used when
  - More than one node performs write in a transaction
  - Explicit 2PC is used
  - More than one node performs write during a single statement
- Only nodes performing writes participate in 2PC
- Design transactions such that they span across as few nodes as possible.

## DBT-2 (sneak peek)

- Like TPC-C
- Early results show 4.3 times scaling with 5 servers
  - More details to come ...



**Thank you**  
ashutosh.bapat@enterprisedb.com