# Poor Man's Parallel Processing

PGConf US 2015 Corey Huinker

#### What is this talk about?

Parallel Processing in Postgres.

### What this talk is really about?

How wonderfully hackable PostgreSQL is.

## Problem: Lack of parallel query in Postgres is hampering adoption.

So, do something about it.

## Aren't there available commercial offerings?

Yes, but that's no fun.

## What about async sharding (PL/Proxy, etc)?

- Have to build your database around the sharding mechanism.
- Nontechnical people laugh when you say "sharding".

#### Common technique: Unix Parallel

- Break up your query into smaller queries.
  - One worker handles A-C, next handles D-F...
- Run them separately, combine the results yourself.
  - o lck.

#### The Goal:

- Something that lets you make something close to an adhoc query.
- Leveraging multiple CPUs on this machine.
- And maybe that other machine too.
- And have the results coalesced into something that can itself be queried (like a table function).
- Without leaving the query.

#### Challenges for general parallelism:

- How should I best break up this big query into smaller ones?
  - With no other information, most systems just do a hash distribution.
- At what point would I overload this machine with worker processes?
- Am I just creating a lot of process/network traffic for myself?
  - Poor distribution means lots of interprocess chatter.

#### PMPP answers *none* of these.

- So why aren't they in PostgreSQL already?
  - Market is littered with problematic parallel halfmeasures.
  - PostgreSQL Hackers want to get it right the first time.
  - Perfect is the enemy of good in this case.
  - Perfect will be nice when we get it (9.5? 9.6?).
  - In the mean time, here's a half-measure that works in limited circumstances if you're careful.

#### What does PMPP look like?

When all of your data is on the same machine, but you want to use multiple CPUs:

And for when you want to query multiple machines:

#### What does PMPP look like? Zoom in.

When all of your data is on the same machine, but you want to use multiple CPUs:

```
function pmpp.distribute( p_row_type anyelement, polymorphic type-spec

A list of SQL statements to be executed.

p_connection text, Any postgres DSN string

p_sql_list text[],

What % of CPUs to allocate.

p_cpu_multiplier float default 1.0 )

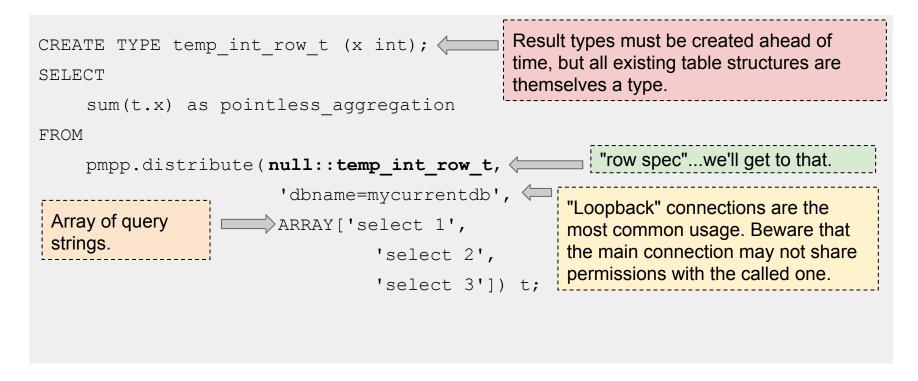
returns setof anyelement
```

And for when you want to query multiple machines:

### What's this null::thingamabob business?

- It's a polymorphic function.
- It gives the shape of the result set that the outer query can expect to receive.
- Is null by convention

#### Example: single machine queries



#### **Example: Query List via Meta-SQL**

```
CREATE TYPE temp int row t (x int);
                                               Just here for an example, you don't have to
                                               redefine it every time.
SELECT
    sum(t.x) as overall rowcount
FROM
    pmpp.distribute(null::temp int row t,
                        'dbname=mycurrentdb',
 Using SQL to generate
                        ARRAY (
                                 SELECT
 SQL is a very powerful
                                       'select count(*) from ' || l.table name
 way to generate worker
 commands. The
                                  FROM
 array() cast helps
                                       partition list 1 )) t;
 visually separate the
 inner and outer queries.
```

#### **Example multi-machine query**

```
SELECT
    sum(t.x) as overall rowcount
FROM
    pmpp.distribute(null::temp int row t,
        '[{"connection":"local dsn", "queries":["SELECT sum(page loads)
    FROM video ads WHERE client = ''CUSTOMER1'' AND ad date >= ''2014-01-
    01''"], "multiplier":"0.5"}, {"connection":"archive dsn", "queries":
    ["SELECT sum(page_loads) FROM video ads WHERE client = ''CUSTOMER1''
    AND ad date < ''2014-01-01''"], "workers": "2"}]'::jsonb) t;
                                 WHA???
```

#### Wait, what was that JSON about?

```
[{"connection": "local dsn", | Each section has connection info, like the local version.
  "queries":[
                                                    We'd normally expect a lot of queries in
     "SELECT sum (page loads) FROM video ads
                                                     at least one of the sections, but this is
     WHERE client = 'CUSTOMER1'
                                                    just an example.
     AND ad date \geq '2014-01-01'''],
                                            We know it has PMPP installed and we want to
  "multiplier":"0.5"},
                                            use AT MOST half the CPUs.
{ "connection": "archive dsn",
  "queries":[
     "SELECT sum (page loads) FROM video ads
                                                       The queries have to all have the
                                                       same shape of result set.
    WHERE client = 'CUSTOMER1'
    AND ad date < '2014-01-01'''],
  "workers":"2"}]
                     Might not have PMPP installed, might not even be real PostgreSQL...
```

## Did you try anything other than polymorphic functions? - Yes: JSON

```
SELECT
                                                                 Re-composition
    sum((t.json data->>'row count')::bigint) as row count
                                                                 acrobatics and
FROM
                                                                 typecasting
    mpp_dist_json( Project name has changed over time
         ARRAY (SELECT
                  'select count(*) as row count from partitions.'
                     || partition name
                                                             Meta-SQL is basically
                FROM
                                                             the same.
                  partition metadata table
                WHERE
                  table name = 'my partitioned table')
                  ) t;
```

It's not the prettiest, and the decompose-recompose overhead increases with the number of columns.

## Did you try anything other than polymorphic functions? - HSTORE

```
SELECT
    sum((t.hstore data->'row count')::bigint)
FROM
    pmpp dist hstore(
        array (SELECT
                  'select count(*) as row count from partitions.'
                    || partition name
               FROM
                 partition metadata table
               WHERE
                 table name = 'my partitioned table')
                 ) t;
```

Basically the same tradeoffs as JSON/JSONB.

#### What's under the hood?

- DBLINK extension
  - dblink\_send\_query() and dblink\_get\_result() async functions
  - This module lacked ability to do polymorphic result sets.
    - So I wrote a patch for that.
    - Ain't hackability great?
- A pg\_attribute query to create table spec
  - FROM dblink get result(x) AS t(col1 int, ...)
  - Query has to be constructed dynamically once, and re-run once per subquery.
  - PL/PGSQL lacks a PREPARE statement
    - Thought about moving to plv8 or C.
  - Will still need this until DBLINK supports polymorphism.

### Under the hood: pg\_attribute query

```
WITH x as (
   select a.attname || ' ' || pg catalog.format type(a.atttypid,
            a.atttypmod) as sql text
   from pg catalog.pg attribute a
   where a.attrelid = pg typeof(p row type)::text::regclass
   and a.attisdropped is false
   and a.attnum > 0
    order by a.attnum )
SELECT format('select * from dblink get result($1) as t(%s)',
                string agg(x.sql text,','))
INTO fetch results query
FROM
      х;
```

Runtime: about 1ms.

#### What's under the hood?

- PL/PGSQL
  - O ONE FOR LOOP
    - really just there to look for failures in initial query distribution.
  - o and one WHILE LOOP
    - looking for queries that have finished, launching new queries as old ones complete, closing down connections
      - pg sleep() with exponential backoff
  - A surprising amount of iteration can be handled in SQL itself.
- temp tables for work queue management, connection management.
  - Wasn't appreciably slower than PL/PGSQL arrays and state variables.
  - Cleaner code, likely very easy to port to C/v8, etc.

### How do you know how many workers to spawn?

By cheating! Hijack the copy command to invoke a command line.

```
create temporary table nproc result (nproc integer);
copy nproc result from program 'nproc'; Sooooo not portable.
format('$$ select greatest(1, (p multiplier * %s)::integer)$$',
          nproc) as nproc sql
from
   nproc result
                   Saves each column of the one-row result set as a
\qset ⇐─
                   same-named variable
create or replace function pmpp.num cpus(p multiplier in float default
1.0) returns integer
                                    Using PSQL vars in SQL definitions.
```

So now you've got an immutable function: ultra-low overhead.

#### How are you using it?

- ETL
  - Partition refresh in place of python & multiprocessing
  - Index Rebuilds
- Deployment scripts
  - Partition creation
- Big-Question queries
  - our data is timeseries, so asking questions across all time can be compute intensive. Partial sums make it more manageable.
- In Development
  - Three-tiered data storage
    - in-memory cache accessed via custom FDW
    - Vertica for recent data
    - Redshift for archive data

#### So many questions!

#### Q. So this would put passwords in the clear, huh?

Yup, anyone with pg\_stat\_activity visibility on the initiating machine could see them.

#### Q. How do you know how many connections are available?

• You don't! (See: Running With Scissors)

### Q. What if the other machine doesn't have pmpp installed? What if the other machine isn't a "real" postgres (Vertica, Redshift)?

Use the num\_workers parameter instead of the multiplier.

#### Q. What's a good multiplier to use?

- 1.0 on AWS EC2s with local SSD drives.
  - Yes, cpu multipliers on Oracle are usually 2x to 4x the number of CPUs.
  - Our queries are very sum-oriented.

#### **Future Direction**

- Put PMPP on PGXN
- 2. CPU detection extension so that we don't rely on nproc existing anymore.
- 3. Get patch to DBLINK accepted into 9.5.
- 4. Become obsolete.