## When PostgreSQL Can't, You Can

- Keith Fiske
- DBA @ OmniTI

http://www.omniti.com http://www.keithf4.com keith@omniti.com @keithf4

## **OmniTI**, Inc

- Full-stack support for high-traffic websites & applications
  - Millions of users
  - Terabytes of data
  - Gilt Groupe, Etsy, Ora.TV, Freelotto
- Surge Conference http://omniti.com/surge
  - Disaster Porn
  - Annually in Sept
- We're hiring!

## **PG Extractor**

- pg\_dump/pg\_restore limitations (-t, -n, -P)
- Filter by schema, table, view, function, type, owner
- Dumps each database object to its own file
- Use regex matching
- Python class (requires python 3)

### **Extensions**

- Introduced in 9.1
- Logically grouped set of database objects
  - CREATE EXTENSION pg\_partman [SCHEMA partman];
- Versioned
  - ALTER EXTENSION pg\_partman UPDATE TO '1.6.1';
  - Update and revert changes predictably.

- Autonomous functions
- Log steps of running function
- Monitors logged functions to ensure they complete
- If/when they fail, where and why

add\_job('job name');

add\_step(job\_id, 'What this step will do'); ... do some stuff... update\_step(step\_id, 'good\_status', 'What this step did successfully');

add\_step(job\_id, 'What this next step will do'); ...do some stuff in a loop... update\_step(step\_id, 'good\_status', 'update every loop iteration to track progress');

add\_step(job\_id, 'One last step'); ... do just a bit more stuff... update\_step(step\_id, 'good\_status', 'Job finished ok');

close\_job(job\_id);

EXCEPTION

WHEN OTHERS THEN

update\_step(step\_id, 'bad\_status', 'Uh..oh...: '||coalesce(SQLERRM,'wat')); fail\_job(job\_id);

#### show\_job('my job name', [int]);

-[ RECORD 3 ]----job id | 10 owner I keith job name | PG JOBMON TEST BAD JOB start time | 2012-09-15 00:55:44.742176-04 end time | 2012-09-15 00:55:44.851514-04 status | CRITICAL pid | 5848 -[ RECORD 4 ]----iob id 19 owner | keith job name | PG JOBMON TEST GOOD JOB start\_time | 2012-09-15 00:55:44.293575-04 end time | 2012-09-15 00:55:44.725483-04 status I OK | 5848 pid

show\_job\_like('I forgot my job's whole name', [int]);

show\_detail(job\_id);
show\_detail('job\_name', [int]);

show\_job\_status('bad\_status', [int]);
show\_job\_status('job name', 'status', [int]);

show\_running([int]);

-[ RECORD 1 ]+-----job id 9 step id | 19 action | Test step 1 start\_time | 2012-09-15 00:55:44.501825-04 end time | 2012-09-15 00:55:44.593389-04 elapsed time | 0.091564 status 1 OK message | Successful Step 1 -[ RECORD 2 ]+----job id |9 step id | 20 action | Test step 2 start time | 2012-09-15 00:55:44.643017-04 end time | 2012-09-15 00:55:44.659336-04 elapsed time | 0.016319 status 1 OK message | Rows affected: 2 -[ RECORD 3 ]+-----19 job id | 21 step id | Test step 3 action start time | 2012-09-15 00:55:44.692518-04 end time | 2012-09-15 00:55:44.7087-04 elapsed time | 0.016182 status I OK message | Successful Step 3

check\_job\_status(interval);

- Make nagios check (command and service configs on my blog)
- Shameless plug http://circonus.com (howto on my blog)

SELECT t.alert\_text || c.alert\_text AS alert\_status FROM jobmon.check\_job\_status() c JOIN jobmon.job\_status\_text t ON c.alert\_code = t.alert\_code;

alert\_status

·····

OK(All jobs run successfully)

alert\_status

CRITICAL(KEITH.SOME\_OTHER\_PROCESS: MISSING - Last run at 2012-09-13 07:17:07.86378-04; KEITH.ANOTHER\_PROCESS: MISSING - Last run at 2012-09-13 07:16:30.169683-04;)

alert\_status

WARNING(KEITH.SOME\_CRITICAL\_PROCESS: RUNNING; )

#### Mimeo

**Per-table Replication Extension** 

"The stencil duplicator or mimeograph machine (often abbreviated to mimeo) is a **low-cost** printing press that works by forcing ink through a stencil onto paper...Mimeographs were a common technology in printing **small quantities**, as in office work, classroom materials, and church bulletins." – Wikipedia



## Mimeo

- Streaming & Log Shipping Replication (omnipitr)
- Per-table replication
  - Snapshot
  - Incremental
  - DML
- Quick replication setup and tear-down
- Installed & run from destination database.
- No superuser required
- Column filter
- Where Condition
- Monitor & Audit Trail w/ PG Jobmon

#### Snapshot

- Whole table replication
- Two tables w/ single view
  - Minimize transactional lock during data migration
  - · Brief exclusive lock to swap view source
- Ideal for small or static tables
- Faster than DML replay if majority of table changes often
- Replicate column changes (new, dropped, type)
- No replication if source data has not changed
- Table
  - Single table, no views
  - Options to handle FK (cascade) and reset sequences
  - Good for dev database

#### Incremental

- Control timestamp column (serial in dev)
- High transaction tables w/ timestamp set every insert
- With primary/unique key, can also support updates
- DST
- Run database in GMT/UTC
- Replication does not run

#### • DML

- Replay Inserts, Updates, Deletes
- Trigger w/ queue table on source
- Doesn't actually replay
  - Queue table of only primary/unique key values for every write
  - Distinct on queue table
  - Delete all matches on destination & re-insert
- Supports multiple destinations

#### • Log Deletes

- Same methods as DML but does not replay deletes
- Common in data warehousing
- Queue table stores entire row if it's deleted
- Destination has special column with timestamp of row's deletion

## **Use Cases**

#### • Table audit

- Trigger to track all changes to audit table w/ audit\_timestamp column
- Use incremental replication on audit table to pull to data warehouse.
- Time-based partitioning on source audit table to easily drop old data with minimal impact on production.
- Database Upgrade
  - Can connect with dblink to any version that supports it
  - Setup replication for larger tables to minimize downtime for pg\_dump upgrade method.

## **PG Partition Manager**

- Current partition management is entirely manual
  - http://www.postgresql.org/docs/current/static/ddl-partitioning.html
- Custom write all tables, triggers, functions, rules, etc.
- Core devs working on getting it built in
  - https://wiki.postgresql.org/wiki/Table\_partitioning
- In the mean time ...

## **Automated Creation**

- Time & Serial Based Partitioning
  - Yearly, Quarterly, Monthly, Weekly, Daily, Hourly, 1/2 hour, 1/4 hour
  - Custom time interval
- Static & Dynamic Triggers
- Pre-creates partitions (customizable how many)
- Manage child table properties from parent
  - Indexes, constraints, defaults, privileges & ownership
- Automatically updates trigger functions as needed.
- Handles object name length limit (63 char)
- Constraint exclusion for non-partition columns

## **Automated Creation**

- Python script to partition existing data
- Commits after each partition created or
- Commit in smaller batches with configured wait
- Partition live, production tables
  - Partitioned 74 mil row table by day (30 days of data)
  - Committed in hourly blocks w/ 5 second wait
  - Streaming slave never fell more than 100 seconds behind
  - 2-3 second lock on parent was only interruption

### **Static Partitioning**

#### Readable functions!

CREATE OR REPLACE FUNCTION partman test.time static table part trig func() **RETURNS** triager LANGUAGE plpasal AS \$function\$ BEGIN IF TG OP = 'INSERT' THEN IF NEW.col3 >= '2013-03-21 00:00:00-04' AND NEW.col3 < '2013-03-22 00:00:00-04' THEN INSERT INTO partman test.time static table p2013 03 21 VALUES (NEW.\*); ELSIF NEW.col3 >= '2013-03-20 00:00:00-04' AND NEW.col3 < '2013-03-21 00:00:00-04' THEN INSERT INTO partman test.time static table p2013 03 20 VALUES (NEW.\*); ELSIF NEW.col3 >= '2013-03-22 00:00:00-04' AND NEW.col3 < '2013-03-23 00:00:00-04' THEN INSERT INTO partman test.time static table p2013 03 22 VALUES (NEW.\*); ELSIF NEW.col3 >= '2013-03-19 00:00:00-04' AND NEW.col3 < '2013-03-20 00:00:00-04' THEN INSERT INTO partman test.time static table p2013 03 19 VALUES (NEW.\*); ELSIF NEW.col3 >= '2013-03-23 00:00:00-04' AND NEW.col3 < '2013-03-24 00:00:00-04' THEN INSERT INTO partman test.time static table p2013 03 23 VALUES (NEW.\*); ELSE **RETURN NEW:** END IF: END IF; **RETURN NULL**; END \$function\$

## **Dynamic Partitioning**

CREATE OR REPLACE FUNCTION partman test.time dynamic table part trig func() **RETURNS** trigger LANGUAGE plpqsql AS \$function\$ DECLARE v count int; v partition name text: v partition timestamp timestamptz; v schemaname text; MAGIC v tablename text: **BEGIN** IF TG OP = 'INSERT' THEN v partition timestamp := date trunc('day', NEW.col3); v partition name := 'partman test.time dynamic table p'|| to char(v partition timestamp, 'YYYY MM DD'); v\_schemaname := split part(v partition name, '.', 1); v tablename := split part(v partition name, '.', 2); SELECT count(\*) INTO v count FROM pg tables WHERE schemaname = v schemaname AND tablename = v tablename; IF v count > 0 THEN EXECUTE 'INSERT INTO '||v partition name||' VALUES(\$1.\*)' USING NEW; ELSE **RETURN NEW:** END IF: END IF: **RETURN NULL;** END \$function\$

### **Automated Destruction**

- Configurable retention policy
  - Time: Drop tables with values older than 3 months
  - Serial: Drop tables with values less than 1000 minus current max
- By default only uninherits
- Can drop old tables or only their indexes
- Dump out tables for archiving
- Python script to undo partitioning

# PgTAP

- Unit testing for PostgreSQL queries, schema & scripting
- Essential for extension development
- ... and your sanity
- http://pgtap.org/

## Links

- https://github.com/omniti-labs/pg\_extractor
- https://github.com/omniti-labs/pg\_jobmon
- https://github.com/omniti-labs/mimeo
- https://github.com/keithf4/pg\_partman

## **Bonus!**

- pg\_bloat\_check.py
- Stol... borrowed query from check\_postgres.pl
- Provide easily readable report on current status of table & index bloat