



PostgreSQL upgrade project

Zdeněk Kotala

Revenue Product Engineer

Sun Microsystems



Agenda

- Overview
- Catalog upgrade
- Storage upgrade
- Others

Overview

Goals

- Minimal downtime
- No extra disk space
- No old version
- Easy to use

Possible design

- Standalone product
 - > Separate binaries which converts database cluster from one version to another.
- PostgreSQL offline upgrade mode
 - > Special mode like bootstrap only on already created cluster.
 - > Data are binary converted.
- PostgreSQL online data conversion
 - > PostgreSQL converts data structure on the fly. Data are converted on background.
 - > PostgreSQL will be able to read old structure.

Possible design

Standalone product

- Advantages
 - > No or minimal impact on core
- Disadvantages
 - > Difficult maintenance – synchronization with core generates a lot of double work
 - > Two code could generate inconsistency
 - > Database is offline during upgrade – downtime depends on database size
 - > Does not fit with PostgreSQL release cycle
 - > No responsibility to implement changes from core

Possible design

PostgreSQL offline upgrade mode

- Advantages
 - > Integrated into core - can reuse server code
 - > No extra application
 - > Middle impact on core (mostly new functions)
 - > **Probably** faster than data export/import
 - > No extra space
- Disadvantage
 - > Database is offline during upgrade - downtime depends on database size

Possible design

PostgreSQL online data conversion

- Advantage
 - > Minimal downtime
 - > Downtime doesn't depend on database size
- Disadvantage
 - > How to convert catalog content and structure
 - > Online data conversion has performance impact depends on implementation (last measure shows 1% performance gap)
 - > Convert on read has lot of corner case problems

...and the winner is

On-line
upgrade

Catalog upgrade

List of affected objects

- Control file
- Flat files
- Directory structure
- Catalog tables
- Configuration

Current solutions

- Pg_migrator or pg_upgrade.sh
 - > Only works for 8.1->8.2
 - > Does not support data layout changes (inet/cidr)
 - > Fast (short downtime)
 - > Problem with tablespaces (keep data on one mount point)
 - > Problem with TOAST tables (TOAST pointer)
 - > Depends on private interfaces

How pg_upgrade.sh works*

- 1) Dump metadata
- 2) Save relation map (relfilenode<->name)
- 3) Export control file data
- 4) Initdb new database cluster
- 5) Freeze database cluster
- 6) Copy CLOG
- 7) Set control data (XID,OID,XLOG ...)
- 8) Create databases, users ...

*Simplified version without tablespaces

How pg_upgrade.sh works (cont.)

- 9) Protect TOAST tables (need to have same relfilenode)
- 10) Create tables, views ...
- 11) Adjust relfilenode for TOAST tables,idx
- 12) Copying and renaming data files
- 13) Done

How upgrade should work

```
pg_ctl -D /var/postgres upgrade
```

How upgrade should work II.

```
check directory /var/postgres ... ok (version 822)
check subdirectories ... ok
creating template1 database in /tmp/pokus/base/1 ... ok
initializing pg_authid ... ok
initializing dependencies ... ok
creating system views ... ok
loading system objects' descriptions ... ok
creating conversions ... ok
creating dictionaries ... ok
setting privileges on built-in objects ... ok
creating information schema ... ok
vacuuming database template1 ... ok
upgrading pg_global database ... ok
upgrading template0 ... ok
upgrading postgres ... ok
upgrading super_db ... ok
```


Control file

- Compatibility verification (BLCKSZ, MAXALIGN, FP format...)
- BLCKSZ, RESEGSIZE, TOAST MAX CHUNK SIZE could be modified during upgrade
- Translate XID, OID, LC_COLLATE, LSN...

Catalogs

- Structure
 - > Use postgres.bki to initialize catalog
 - > Keep old data files for data transfer
- Contents
 - > User metadata will be transferred and converted to the new structure
 - > Strict rules for using same OID for modified or new object
 - > Some kind of changes is not allowed (e.g. binary format change must invoke new data type - new OID)

Configuration files

- postgresql.conf
 - > New GUC variable will contain default value
 - > Obsolete GUC variable will be ignored - warning in log file
 - > Out of range values will be set to default
 - > Problem is with different meaning of values
- pg_hba.conf, pg_ident.conf
 - > Depends on kind of change ...

Storage upgrade

Page Layout Structures

BLCKSZ

PageHeaderData

TOAST_MAX_CHUNK_SIZE

ItemIdData

*MaxItemSize

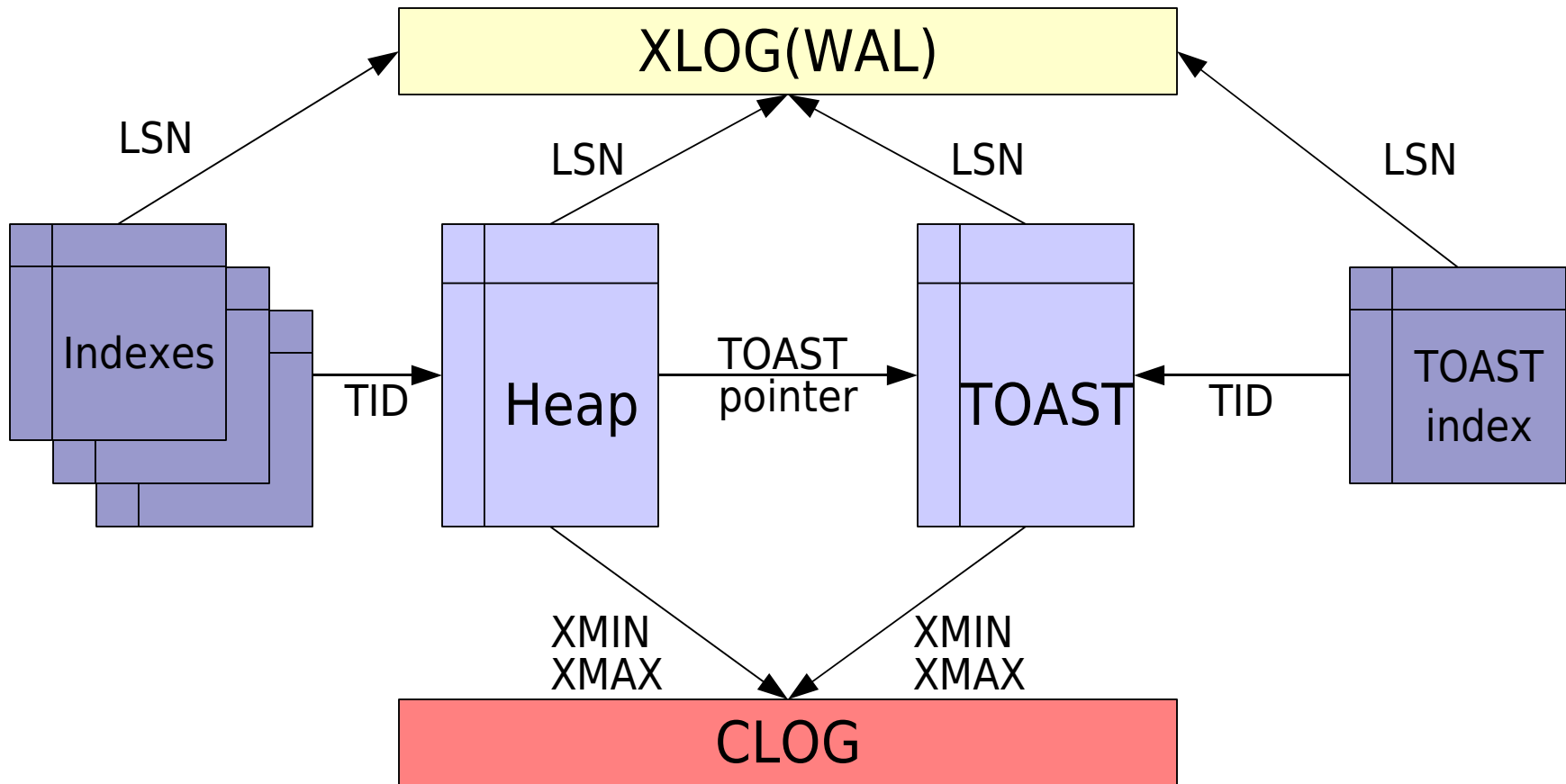
IndexTupleData

*OpaqueData

HeapTupleHeaderData

varatt*

Storage dependency graph



Issues to solve

- Any storage relation **MUST NOT** be break
- Data have to be converted or they have to contain version info
 - > Page has page version in a header
 - > HeapTupleHeader and IndexTupleHeader does not have this information
 - > Big problem with toasted arrays and composite datatypes. There is not clue what kind of data are stored in the chunks. Page cannot be converted when page is shared with more chunked datums.

Storage upgrade methods

- On line
 - > Read only mode
 - > Read All, Write New
 - > On fly page layout conversion
- Off line
 - > Dump/Restore
 - Need extra space
 - > Heap conversion+Retoasting+Reindexing
 - Needs lot of extra steps and probably it is not speedup

Storage upgrade methods

Offline

- TOAST table upgrade
 - > Upgrade TOAST table pages and chunk header
 - > Reindex TOAST table index
- Heap table upgrade
 - > Convert tuples include convert related toasted values
- Reindex table

Storage upgrade methods

Read Only Mode

- Need to learn PostgreSQL to work with old data structures
- Add extra code which could slow down general performance
- Easy return back to prior version
- Problem with catalog

Storage upgrade methods

Read All, Write New

- Based on Read Only Mode
- New and modified data are written in new format
- Storage/access layer returns back HeapTupleData with version information and executor converts data to the new format when it is necessary
- General performance has ~1% drop

Storage upgrade methods

Read Old Write New - example

```
#define SizeOfPageHeaderData(page) \
    (PageGetPageLayoutVersion(page) == 4 ? \
        (offsetof(PageHeaderData_04, pd_linp[0])) :\
        (offsetof(PageHeaderData_03, pd_linp[0])))
```

```
typedef struct HeapTupleData
{
    uint32    t_len;           /* length of *t_data */
    ItemPointerData t_self;   /* SelfItemPointer */
    Oid       t_tableOid;     /* table the tuple came from */
    uint16    t_version;      /* page layout version */
    HeapTupleHeader t_data; /* -> tuple header and data */
} HeapTupleData;
```

Storage upgrade methods

Online Page Layout Conversion

- Possible only when converted data fits on same page and toasted data structures are not affected.
- “Not possible” between layout version 3 and 4 (8.2->8.3).
 - > Pageheader has been extended to 24 bytes.
 - > Index tuples does not fit on a page, different toast chunk size and heap tuples does not fit on machines with MAXALIGN=8 (e.g. SPARC, 64bit x86)
- WAL generates a lot of full page writes.

Storage upgrade methods

Online Page Layout Conversion - example

- Converter hook in ReadBuffer_common

```

{
    smgrread(reln->rd_smgr, blockNum, (char *) bufBlock);
    /* Page Layout Converter hook. We assume
       that page version is on same place. */
    if( plc_hook && PageGetPageLayoutVersion(reln,bufBlock)
        != PG_PAGE_LAYOUT_VERSION )
    {
        plc_hook((char *)bufBlock);
        bufHdr->flags |= (BM_DIRTY | BM_JUST_DIRTIED);
        log_newpage(&reln->rd_node, blockNum ,bufBlock);
    }
}

```

Storage upgrade methods

Inner heap tuples reorganization

- Similar to page layout conversion, but tuple which does not fit on the page have to be moved to a new page
- Need handle inter page transfer
 - > Reindex
 - > Introduce inter page redirection pointer
 - > Mark tuple dead and insert new tuple to a free page in the new format
- Requires WAL logging
- Still problem with toasted values

Storage upgrade methods

Retoasting

- Needed when `TOAST_MAX_CHUNK_SIZE` has been changed or when arrays and composite datatype need to be converted
- More possible solutions
 - > Add `TOAST_MAX_CHUNK_SIZE` to `pg_class`
 - > Adjust `toast_fetch_datum()` accept different chunk size
- Need to convert TOAST table index
- No clue what datatype is stored in chunks

Storage upgrade methods

Reindexing

- Reindexing is necessary every time when
 - > tid of any tuple has been changed
 - > index structure has been changed
 - > index tuples does not fit on a new page layout
 - AM interface could have convert function for each index type
- Reindex could be performed on the running system

Write Ahead Log (WAL/XLOG)

- CHECKPOINT is last operation on shutdown. All changes are applied and WAL files can be dropped.
- Needs to keep XLOG pointer to protect correct recovery (LSN dependency on WAL)

Commit log (CLOG)

- Array of transactions status
- No changes for long time - stable
- Some upgrade methods could produce a frozen database, afterwards CLOG files could be removed

Other

Stored procedures

- Changes in PL languages
 - > All changes are usually backward compatible
 - > Possible to add language version into catalog and delivery more *.so
 - > Problem with procedures written in C

Tsearch2

- Any change in FTS configuration or dictionary implies regeneration of affected tsvector fields. Unfortunately, there is not relation between tsvector and original source.

Version 8.4 breakers

- New HASH index implementation
- Integer datetime
- CRC
- New tuple alignment

Conclusion

Conclusion

- Only on-line in-place upgrade is acceptable
- Still not clear what method is better. However many modifications are same for both approaches.
- Needs lot of cleanup changes which requires acceptance from community and committers
- I still hope that it will be integrated in to 8.4

References

<http://pgfoundry.org/projects/pg-migrator/>

<http://src.opensolaris.org/source/xref/sfw/usr/src/cmd/postgres/postgresql-upgrade/>

http://wiki.postgresql.org/wiki/In-place_upgrade



PostgreSQL upgrade project

Zdeněk Kotala

zdenek.kotala@sun.com

