

Optimiser la table postgresql `ejabberd.archive`

22/10/2022 12:13 - Anonyme

Statut:	Nouveau	Début:	22/10/2022
Priorité:	Normale	Echéance:	
Assigné à:	pitchum .	% réalisé:	0%
Catégorie:		Temps estimé:	0.00 heure
Version cible:	Backlog		
Difficulté:	5 Difficile		

Description

Le problème

Cette table `ejabberd.archive` conserve l'historique de tous les messages (dans XMPP ça s'appelle MAM - Message Archive Management).

Elle est donc très volumineuse. Et elle est très souvent sollicitée, aussi bien en lecture qu'en écriture.

Des "SELECT" ou "SELECT COUNT" sur cette table durent souvent plus de 30 secondes, parfois plusieurs minutes. Ça doit probablement se ressentir côté client (mais je n'en suis pas certain, c'est peut-être habilement masqué dans les applis).

Une piste de solution

Une vraie solution nécessiterait de patcher le code ejabberd. Par exemple pour maintenant des compteurs de messages pour éviter d'avoir à faire des "SELECT COUNT" violents.

Mais comme je ne suis pas prêt à apprendre à code en erlang ni en elixir, je cherche ce que je peux faire côté postgresql.

Une amélioration possible serait de "saucissonner" la table (ça s'appelle [Table Partitioning](#) dans la doc postgresql).

Ce mécanisme se prête à priori très bien à cette table pour laquelle les écritures sont de type "append only" et qui dispose d'un champ d'horodatage.

L'idée serait de regrouper les messages par trimestres par exemple. Ça devrait rendre beaucoup plus rapides certaines requêtes SELECT basés sur l'horodatage.

Début d'analyse

Avec le module pg_stat_statements activé depuis un bon moment maintenant, on peut avoir un aperçu des requêtes les plus fréquentes et leurs durées moyennes d'exécution.

Je filtre les requêtes qui correspondent aux backups et à la métrologie. Leurs durées d'exécution nous importent peu.

```
postgres=# SELECT
  substring(query, 1, 120) as short_query,
  ROUND(total_exec_time :: numeric, 2) AS total_time,
  calls,
  ROUND(mean_exec_time :: numeric, 2) AS mean,
  ROUND (
    ( 100 * total_exec_time / SUM (total_exec_time :: numeric) OVER () ) :: numeric,
    2
  ) AS percentage_overall
FROM pg_stat_statements
WHERE calls > 10
  AND query NOT LIKE 'COPY %'
  AND query NOT LIKE '% where peer not like %'
  AND query NOT LIKE 'alter database %'
ORDER BY total_time DESC
LIMIT 20;
```

	short_query	total_time	calls	mean	percentage_overall
-----+-----+-----+-----+-----					
SELECT COUNT(*) FROM archive WHERE username=\$1 and server_host=\$2		30119788.74	357780	84.19	60.09

```

SELECT timestamp, xml, peer, kind, nick FROM (SELECT timestamp, xml, peer, kind, nick FROM archive WHERE username=$1 and server_host=$2 AND timestamp >= $3 ORDER BY timestamp)
INSERT INTO archive(username, server_host, timestamp, peer, bare_peer, xml, txt, kind, nick) VALUES ($1, $2, $3, $4, $5, $6, $7, $8, $9)
SELECT timestamp, xml, peer, kind, nick FROM archive WHERE username=$1 and server_host=$2 AND timestamp >= $3 ORDER BY timestamp
select feature from caps_features where node=$1 and subnode=$2
insert into caps_features(node, subnode, feature) values ($1, $2, $3)
INSERT INTO spool(username, server_host, xml) VALUES ($1, $2, $3)
select name, val from pubsub_node_option where nodeid=$1
SELECT COUNT(*) FROM archive WHERE username=$1 and server_host=$2 and bare_peer=$3 and timestamp >= $4 and timestamp <= $5
select password, serverkey, salt, iterationcount from users where username=$1 and server_host=$2
select itemid, publisher, creation, modification, payload from pubsub_item where nodeid=$1
WITH upsert AS (UPDATE last SET seconds=$1, state=$2 WHERE username=$3 AND server_host=$4 RETURNING *) INSERT INTO last(username, server_host, seconds, state) VALUES ($3, $4, $1, $2)
select node, parent, plugin, nodeid from pubsub_node where host=$1
SELECT COUNT(*) FROM archive WHERE username=$1 and server_host=$2 and timestamp <= $3
delete from spool where username=$1 and server_host=$2
select node, plugin, i.nodeid, affiliation from pubsub_state i, pubsub_node n where i.nodeid = n.nodeid and jid=$1 and h=$2
select username, jid, nick, subscription, ask, askmessage, server, subscribe, type from rosterusers where username=$1 and server_host=$2
WITH upsert AS (UPDATE pubsub_item SET publisher=$1, modification=$2, payload=$3 WHERE nodeid=$4 AND itemid=$5 RETURNING *) INSERT INTO pubsub_item(nodeid, itemid, publisher, modification, payload) VALUES ($4, $5, $1, $2, $3)
SELECT timestamp, xml, peer, kind, nick FROM (SELECT timestamp, xml, peer, kind, nick FROM archive WHERE username=$1 and server_host=$2 AND timestamp >= $3 ORDER BY timestamp)
select itemid, publisher, creation, modification, payload from pubsub_item where nodeid=$1 order by creation asc
(20 lignes)

```

Les cas d'usage typiques que je pense avoir identifiés :

- une appli se déconnecte/reconnecte, même brièvement, elle vérifie immédiatement si son historique de discussion est à jour => "SELECT COUNT FROM archive"
- pareil, mais en pire, si un utilisateur configure son compte XMPP déjà existant sur un nouvel appareil, celui-ci va potentiellement chercher à récupérer tout l'historique de ce compte
- un utilisateur rallume un appareil inactif depuis longtemps, l'appli va interroger MAM pour récupérer tous les messages qui ont été archivés depuis tout ce temps : "SELECT FROM archive WHERE timestamp >="
- un utilisateur reçoit ou envoie un message => "INSERT INTO archive ..."
- un utilisateur rejoint un salon existant => "SELECT COUNT ... FROM archive WHERE ..."
- certaines applis ne récupèrent l'historique que au fil de l'eau, lorsque l'utilisateur scrolle vers le haut : "SELECT FROM archive WHERE timestamp >= AND timestamp <= ..."

Plan d'action

Yapuka ;)