

UNIVERSITATEA ALEXANDRU IOAN CUZA IAȘI

FACULTATEA DE INFORMATICĂ



LUCRARE DE DISERTAȚIE

A study on point in time recovery on relational databases

propusă de

Ștefan Stan

Sesiunea: *Iulie, 2017*

Coordonator științific

Lect. dr. Cristian Frăsinaru

UNIVERSITATEA ALEXANDRU IOAN CUZA IAȘI

FACULTATEA DE INFORMATICĂ

A study on point in time recovery on relational databases

Ștefan Stan

Sesiunea: *Iulie, 2017*

Coordonator științific

Lect. dr. Cristian Frăsinaru

Supervizor companie

DevOps Engineer Gabi Ichim, TSS-Yonder

DECLARAȚIE PRIVIND ORIGINALITATEA ȘI RESPECTAREA DREPTURILOR DE AUTOR

Prin prezenta declar că lucrarea de disertație cu titlul "A study on point in time recovery on relational databases" este scrisă de mine și nu a mai fost prezentată niciodată la o altă facultate sau instituție de învățământ superior din țară sau străinătate. De asemenea, declar că toate sursele utilizate, inclusiv cele preluate de pe Internet, sunt indicate în lucrare, cu respectarea regulilor de evitare a plagiatului:

- toate fragmentele de text reproduse exact, chiar și în traducere proprie din altă limbă, sunt scrise între ghilimele și dețin referința precisă a sursei;
- reformularea în cuvinte proprii a textelor scrise de către alți autori deține referința precisă;
- codul sursă, imagini etc. preluate din proiecte open-source sau alte surse sunt utilizate cu respectarea drepturilor de autor și dețin referințe precise;
- rezumarea ideilor altor autori precizează referința precisă la textul original.

Iași,

Absolvent *Ștefan Stan*

DECLARAȚIE DE CONSIMȚĂMÂNT

Prin prezenta declar că sunt de acord ca lucrarea de disertație cu titlul "A study on point in time recovery on relational databases", codul sursă al programelor și celelalte conținuturi (grafice, multimedia, date de test etc.) care însoțesc această lucrare să fie utilizate în cadrul Facultății de Informatică.

De asemenea, sunt de acord ca Facultatea de Informatică de la Universitatea Alexandru Ioan Cuza Iași să utilizeze, modifice, reproducă și să distribuie în scopuri necomerciale programele-calculator, format executabil și sursă, realizate de mine în cadrul prezentei lucrări de disertație.

Iași,

Absolvent *Ștefan Stan*

Cuprins

Introducere	6
Contribuții.....	7
1 Point in time recovery.....	8
1.1 Setarea arhivării fișierelor WAL.....	9
1.2 Crearea unui basebackup	11
1.3 Restabilirea db. folosind Continuous Archive Backup	12
2 Logical Volume Manager (LVM).....	13
2.1 Arhitectura LVM.....	13
2.2 Funcționalități LVM.....	13
2.3 Crearea Logical Volumes.....	14
2.3.1 Pași pentru a crea un Physical Volume.....	14
2.3.2 Pași pentru a crea un Volume Group.....	15
2.3.3 Pași pentru a crea Logical Volumes	15
2.3.4 Modificarea mărимii unui Logical Volume.....	15
2.4 Explicații despre convenția de numire a LVM.....	17
3 Database Point in time recovery Manager	18
3.1 Arhitectura sistemului de fișiere.....	18
4 DevOps as a Service	20
4.1 Arhitectura Sistemului Software	20
4.2 Scripturile UNIX.....	21
4.3 Server-ul: RESTful API.....	23
4.4 Interfața grafică cu utilizatorul	25
4.5 Setarea aplicației în Amazon Web Services.....	26
5 Ghid de utilizare a aplicației	27
6 Direcții de dezvoltare viitoare	28
Concluzii.....	29
Bibliografie.....	30

Introducere

Având în vedere contextul unui dezvoltator software, am observat în timpul experienței faptul că pentru echipa efectivă de dezvoltare era destul de dificil de întreținut baza de date (incluzând deasemenea și operațiile de pornire/oprire, mai ales pentru juniori acest lucru putea pune probleme dacă baza de date se afla instalată pe o platforma UNIX). În cazul în care se dorea aplicarea unei strategii de point in time recovery pentru restabilirea bazei de date la o stare existentă la un moment dat, era nevoie de intervenția echipei de devOps, care se ocupa în mod expres de acest tip de operații, având cunoștințe avansate despre subiect.

Astfel că, am observat necesitatea existenței unei aplicații care să ușureze operațiile asupra bazei de date, începând cu pornirea/oprirea/verificarea stării, până la întreaga operație de backup și point in time recovery (întoarcere la o stare existentă la un anumit moment dat). Această aplicație trebuia să fie intuitivă, ușor de folosit chiar și pentru cineva care nu avea cunoștințe de devOps, nefiind necesar accesul efectiv pe mașina de UNIX unde se afla clusterul de baze de date. Este de la sine înțeles că aplicația ar trebui să detecteze toate bazele de date din clusterul de pe mașina pe care este încărcată și ar trebui să permită operațiile prezentate mai sus pentru fiecare instanță în parte.

Impactele unei astfel de aplicații sunt atât de natură tehnică, le oferă dezvoltatorilor mai mult control asupra sistemului necesitând mai puține cunoștințe, precum și de business, nu mai trebuie angajată o echipă dedicată de devOps și prin urmare se reduc costurile.

Contribuții

În crearea aplicației am văzut o necesitate existentă întâlnită la dezvoltatorii de software și am încercat să o rezolv. Procesul general de restabilire a unei baze de date este descris în *Capitolul 1*.

În *Capitolul 2* se prezintă într-o manieră detaliată o metodă de organizare a sistemului de fișiere UNIX, urmând ca în *Capitolul 3* să fie discutată arhitectura aleasă pentru dispunerea clusterului de baze de date pe sistemul de fișiere, precum și avantajele aduse de aceasta.

Soluția software este trecută în revistă în *Capitolul 4*. Astfel că, se prezintă arhitectura aplicației, urmată de fiecare componentă ce intră în alcătuirea ei. Tot în acest capitol este descrisă și setarea aplicației și a mediului de lucru în Amazon Web Services.

Capitolul 6 oferă un ghid detaliat al aplicației, urmând ca în *Capitolul 7* să fie discutate direcțiile de dezvoltare viitoare și idei de îmbunătățire ale platformei dezvoltate.

Point in time recovery

”Point in time recovery” în contextul informaticii se referă la funcționalitatea unui sistem folosită de un administrator pentru a restabili sau recupera un set de date existent la un moment dat.

PostgreSQL menține în fiecare moment un *write ahead log* WAL în subdirectorul `pg_xlog` al directorului de bază unde avem clusterul de baze de date. Acest log reține fiecare schimbare făcută în fișierele bazei de date. Există, în principiu, pentru momentele în care sistemul cedează. Dacă acest lucru se întâmplă, datele pot fi recuperate prin o executare a secvenței instrucțiunilor din log de la momentul ultimei salvări până înainte de momentul critic. Totuși, existența acestui log face posibilă o a treia strategie de restaurare a bazei de date: putem combina un backup la nivel de sistem de fișiere a bazei de date cu un backup al fișierelor WAL. Dacă avem nevoie să recuperăm datele vom restaura backup făcut la nivel de sistem de fișiere și apoi vom executa o secvență de fișiere WAL pentru a aduce baza de date la o versiune stabilă. Această abordare este complexă și dificil de setat, dar are multe beneficii:

- Nu avem nevoie de o stare consistentă a backup la nivel de sistem de fișiere. Orice inconsistentă va fi remediată prin executarea secvenței de fișiere WAL (nu este cine știe ce diferit de ceea ce se întâmplă în momentul în care sistemul își revine după ce a cedat). Astfel că, avem nevoie de un utilitar capabil să creeze salvări pe disc, tar zip etc.
- De vreme ce putem folosi o secvență nedefinită de fișiere WAL, backup continuu se poate realiza pur și simplu arhivând fișierele WAL. Această tehnică este foarte folositoare în cazul bazelor de date mari, unde nu este avantajos să facem backup complet în mod frecvent.
- Nu este necesar să executăm toată secvență de fișiere WAL. Putem să ne oprim la orice moment din timp în care baza de date era consistentă. Astfel că, această tehnică suportă *point in time recovery*. Este posibil să restaurăm baza de date la orice moment existent dacă avem fișiere WAL.
- Dacă oferim încontinuu fișiere WAL unei mașini noi ce a fost încărcată de la aceeași salvare, putem avea un sistem aflat în hibernare pe care îl putem aduce în orice moment în funcțiune având setul de date aflat în producție.

Pentru a recupera date cu succes folosind strategia de arhivare continuuă este nevoie de o secvență continuuă de fișiere WAL arhivate ce există cel puțin de la momentul în care am creat acel basebackup (o stare a bazei de date aflată pe disc, un checkpoint). Pentru a începe e nevoie să setăm și să testăm procedura de salvare a fișierelor WAL, înainte de a face primul basebackup. În cele ce urmează vom discuta mecanismul de arhivare a fișierelor WAL.

1.1 Setarea arhivării fișierelor WAL

Într-un mod abstract, un sistem PostgreSQL produce o secvență nedefinită de înregistrări WAL. Sistemul împarte această secvență în fișiere WAL, care sunt în mod normal de 16MB fiecare (deși dimensiunea segmentului poate fi alterată când facem build la PostgreSQL). Fișierele segment au numele într-un format numeric care reflectă poziția lor în secvența abstractă WAL. Când nu folosim arhivarea fișierelor WAL sistemul creează câteva fișiere pe care le reciclează redenumindu-le pe cele vechi și marcându-le pentru ștergere.

Când arhivăm datele WAL trebuie să prindem conținutul fiecărui segment de îndată ce devine complet și să salvăm datele undeva până ca segmentul devine reciclat și refolosit. În funcție de aplicație și hardware folosit, pot exista mai mult moduri să ”salvăm datele undeva”: putem copia segmentul într-o locație pe un disc extern, le putem trimite în cloud pe un cont de Amazon S3, Google Drive... le putem trimite prin rețea oriunde etc. PostgreSQL încearcă să nu facă supoziții despre modul de arhivare a datelor pentru a facilita munca administratorului. În schimb, PostgreSQL lasă administratorul să ofere o comandă shell ce se va executa pentru a copia orice segment complet oriunde trebuie să meargă. Această comandă poate fi la fel de simplă ca un **cp**, sau se poate invoca un script complex de linux.

Pentru a activa arhivarea WAL trebuie să se parcurgă o serie de pași:

- să se seteze **wal_level** ca fiind **archive** sau mai mare. (wal_level determină câtă informație se scrie în fiecare WAL. Implicit e setat ca fiind **minimal**, ceea ce înseamnă că se scrie doar informația necesară pentru a se restaura după o problemă sau după o închidere imediată; **archive** adaugă informații necesare pentru arhivarea WAL; **hot_standby** adaugă informații necesare pentru a rula interogări în mod citire pe un server aflat în pauză; **logical** adaugă informații necesare pentru a suporta decodare logică)
- să se specifice o comandă shell în parametrul **archive_command**. În comanda de arhivare, parametrul %p este înlocuit cu calea către fișierul

WAL, în timp ce parametrul %f este înlocuit cu numele fișierului. Calea este relativă la fișierul curent (fișierul în care se află clusterul de baza de date). Cea mai simplă comandă este de genul:

```
archive_command = 'test ! -f /mnt/server/archivedir/%f && cp %p /mnt/server/archivedir/%f'
```

Comanda va copia segmentul WAL în directorul /mnt/server/archvedir.

După ce s-au înlocuit parametrii %p si %f cmd. ce fa fi exec. va arata astfel:

```
test ! -f /mnt/server/archivedir/00000001000000A9000000065 && cp pg-  
xlog/00000001000000A9000000065 /mnt/server/archivedir/00000001000000A9000000065
```

O comandă similară va fi executată pentru fiecare segment nou ce fa fi generat. Comanda de arhivare va fi executată sub aceleași utilizator ca și serverul PostgreSQL. De vreme ce seria de fișiere WAL conține efectiv datele din baza de date, va fi nevoie ca datele arhivate să fie protejate de către accesul neautorizat. Putem să nu oferim drept de citire pentru grup și pentru ceilalți.

Comanda de arhivare este important să returneze statusul 0. PostgreSQL, cand primește 0 va crede că fișierul a fost arhivat cu succes și îl va șterge/ marca pentru reciclare. Însă, un status diferit de 0 va spune PostgreSQL ca arhivarea a eșuat și acesta va păstra fișierul reîncercând ulterior.

Comanda de arhivare ar trebui construită în așa fel încât să refuze să suprascrie fișiere pre-existente. Această funcționalitate este o metodă de siguranță prin care se asigură integritatea datelor în cazul unei erori a administratorului (trimite rezultatul de pe doua servere în același director).

Este de preferat să se testeze comanda propusă pentru arhivare, astfel că aceasta sa nu suprascrie fișiere și să returneze status diferit de 0 cand se încearcă acest lucru. Comanda oferită ca exemplu mai sus asigură acest lucru prin pasul de testare inițial. Unix deține pe anumite platforme opțiuni precum **-i** ce pot fi folosite pentru a realiza același lucru, dar nu trebuie să ne bazam pe ele ci să verificăm codul returnat de fiecare dată.

Aceste setări vor fi modificate în fișierul **postgresql.conf**

1.2 Crearea unui basebackup

Un basebackup este o salvare pe disc la un moment dat a bazei de date, de la care se va porni procesul de executare a fișierelor WAL.

Cel mai ușor mod de a crea un basebackup este folosirea utilitarei **pg_basebackup**. Se poate crea un basebackup fie ca și fișiere normale pe disc sau ca și arhivă(tar). Nu este neapărat nevoie să fim îngrijorați de timpul necesar creării basebackup. Totuși, dacă serverul rulează cu opțiunea **full_page_writes** dezactivată s-ar putea observa o scădere în performanță, de vreme ce procedura de basebackup rulează automat cu opțiunea **full_page_writes** activată în mod forțat.

Pentru a ne folosi de basebackup trebuie să reținem toate fișierele WAL din timpul creării basebackup și după. Pentru a ne ajuta în acest proces, procedura de basebackup creează un fișier backup_history ce este mutat imediat în locația fișierelor WAL. Fișierul este denumit după primul fișier care este necesar pentru recovery folosind acel basebackup. Odată ce procesul de basebackup s-a terminat cu succes nu mai este nevoie de fișierele WAL de dinainte de crearea basebackup, de vreme ce avem o stare stabilă a bazei de date la momentul creării basebackup, de la care se va porni un proces ulterior de executare a fișierelor WAL pentru revenirea în caz de dezastru a bazei de date.

De vreme ce trebuie să salvăm fișierele WAL arhivate până la momentul ultimului basebackup înseamnă că intervalul între care se creează basebackups trebuie ales în funcție de cât spațiu de stocare vrem să folosim pentru a salva fișierele WAL. De asemenea, este de considerat și timpul pe care vrem să îl alocăm unui eventual proces de restabilire a bazei de date. Dacă este necesar acest lucru va trebui să executăm fiecare fișier WAL și ar putea dura destul de mult în caz că ultimul basebackup a fost făcut acum mult timp.

1.3 Restabilirea db. folosind Continuous Archive Backup

Să presupunem că ce e mai rau s-a întâmplat și e nevoie să restabilim baza de date. Procedura este următoarea:

1. Se oprește serverul, dacă rulează.
2. Se copie conținutul fișierul **pg_xlog** ce poate conține fișiere log cu motivele pentru care sistemul a ajuns într-o stare inconsistentă, precum și segmentele WAL nearhivate (ce nu au atins pragul necesar (implicit 16MB) pentru a fi arhivate)
3. Se șterge orice fișier din locația clusterului de baze de date.
4. Se copie datele ultimului basebackup în locația clusterului de baze de date. Trebuie verificate drepturile de acces la fișierele/directoarele copiate.
5. Se șterge conținutul directorului **pg_xlog** copiat din basebackup, deoarece nu este de actualitate. Se copie fișierele din directorul **pg_xlog** salvat anterior.
6. Se creează un fișier **recovery.conf**. De asemenea se va modifica și fișierul **pg_hba.conf** pentru a nu da voie utilizatorilor obișnuiți să se conecteze până când procesul de recuperare nu se termină.
7. Se pornește server-ul. Acesta va intra automat în mod de recuperare și va începe să execute fișierele WAL, dacă există. Dacă procesul eșuează din motive externe, se poate restarta serverul și va continua. În momentul în care procesul se termina se va redenumi fișierul **recovery.conf** în **recovery.done** pentru a preveni porniri ulterioare ale procesului de recuperare. Baza de date va funcționa normal de acum.
8. Se inspectează baza de date pentru a se vedea dacă starea este una consistentă și dorită.

Procesul se bazează pe setarea unui fișier de configurare **recovery.conf** ce specifică cum să se execute procesul și până când anume. Singurul lucru neaparat necesar în fișierul de restaurare este **recovery_command** care specifică ce comandă shell se va executa pentru a recupera fișierele WAL. Aceasta este construită precum **archive_command** prezentată anterior.

De exemplu: **restore_command = 'cp /mnt/server/archivedir/%f %p'**

Logical Volume Manager (LVM)

În zilele de început ale serverelor Linux se făcea managementul spațiului de stocare prin crearea de partiții pe disc. Chiar dacă această abordare funcționează, ea prezintă anumite dezavantaje, iar cea mai importantă este aceea că discurile sunt inflexibile. Acesta este motivul pentru care a fost introdus Logical Volume Manager. Chiar dacă nu este posibil să schimbăm în mod dinamic mărimea unei partiții care rămâne fără spațiu, acest lucru este rezolvat lucrând cu LVM. LVM oferă multe avantaje ce vor fi prezentate în cele ce urmează.

2.1 Arhitectura LVM

În arhitectura LVM se pot distinge mai multe layer-uri. Primul nivel lucrează direct cu componentele de stocare. Acesta este creat într-un mod abstract și poate folosi orice medii de stocare.

Mediile de stocare trebuie marcate ca "physical volumes", ceea ce le face să poată fi folosite de către LVM. O componentă de stocare marcată ca "physical volume" poate fi adăugată unui "volume group", layer abstract peste toate mediile de stocare. Abstractizarea constă în faptul că acesta nu este ceva fix, poate fi redimensionat când este necesar, ceea ce face posibil să se adauge mai mult spațiu la nivelul "volume group" atunci când "logical volumes" ce fac parte din el rămân fără spațiu. Concret, dacă se rămâne fără spațiu pe un logical volume se ia din spațiul liber de la nivelul volume group-ului din care face parte. În cazul în care nici la acel nivel nu mai există spațiu liber se adaugă un disc nou în volume group.

Intuitiv, deasupra nivelului de volume group se află logical volumes. Acestea nu acționează ca medii de stocare directe, ele își iau spațiul necesar din discurile aflate în volume group. Astfel că, un logical volume poate avea spațiul de stocare distribuit pe mai multe physical volumes (medii de stocare).

2.2 Funcționalități LVM

Există multe motive pentru care LVM este foarte util. Cel mai important motiv este acela că oferă flexibilitate pentru managementul spațiului de stocare. Partițiile nu mai sunt limitate de către spațiul de stocare al discurilor. De asemenea, este posibil să se reducă mărimea unui logical volume, dar doar dacă sistemul de fișiere ce a fost creat suportă redimensionare. De exemplu, sistemul de fișiere Ext4 suportă redimensionare în timp ce XFS, care este sistemul de fișiere implicit de pe

RHEL (Red Hat Enterprise Linux) nu suporta.

Un alt motiv pentru care administratorii de sistem folosesc LVM este posibilitatea de a crea salvari. O salvare tine minte starea curenta a logical volume si poate fi folosita pentru il a restaura daca este nevoie.

Al treilea avantaj important este acela ca avem posibilitatea de a schimba usor elementele hardware stricate. Daca un hard disk se strica, datele pot fi mutate in cadrul volume group prin comanda pvmove si un nou harddisk se poate adauga dinamic fara a avea sistemul cazut.

2.3 Crearea Logical Volumes

Crearea Logical Volumes presupune crearea a trei layere in arhitectura LVM. In primul rand trebuie create PV(physical volumes), apoi urmeaza sa se creeze VG(volume group) la care se adauga physical volumes create anterior. In ultimul rand, trebuie create LV(logical volumes) pe care le vom folosi ca si "partitii" efective.

Inainte de a folosi utilitarele LVM trebuie sa marcam partițiile hardware ca fiind de tipul LVM. Daca se foloseste utilitara fdisk impreuna cu un disc de tip MBR, atunci partitia trebuie pusa de tipul 8e. Altfel, daca se foloseste gdisk impreuna cu un disc de tip GUID, atunci partitia este de tip 8300.

2.3.1 Pasi pentru a crea un Physical Volume

Avem nevoie de un harddisk care are spatiu liber (nepartitionat). Pentru scopul explicarii se va folosi un disk /dev/vdb pentru a crea partitia. Fiecare disc are un nume diferit astfel ca pasii vor trebui ajustati pentru fiecare disc.

- Se deschide un shell si se executa **fdisk /dev/vdb**
- Se apasa **n** pentru a crea partitia. Se selecteaza **p** pentru a o marca partitie primara si se accepta numerotarea implicita.
- Se apasa **Enter** de doua ori si astfel se creeza o partitie folosind tot spatiul de pe discul /dev/vdb
- Fiind intors la prompter-ul fdisk, se apasa **t** pentru a schimba tipul partitiei create. Fiind o singura partitie ce foloseste intreg spatiul aceasta va fi selectata implicit, altfel ar putea intreba care pentru care partitie sa se schimbe tipul.
- Utilitara fdisk intreaba ce tip va avea partitia aleaa. Scrie 8e. Apasa **w** pentru a scrie schimbarile efectuate in acesti pasi. Daca se primeste un

mesaj conform caruia tabela de partitii nu a putut fi actualizata trebuie dat un restart.

Acum ca partitia a fost creata trebuie marcata ca LVM physical volume. Pentru asta se va executa **pvcreate /dev/vdb1**. Ar trebui primit urmatorul mesaj: Physical volume “/dev/vdb1” successfully created. Pentru a verifica existenta physical volume nou creat putem executa in shell comanda pvs ce listeaza physical volume existente.

2.3.2 Pasi pentru a crea un Volume Group

Acum ca am creat un physical volume trebuie sa il asignam unui volume group. Este posibil sa adauga physical volumes la un volume group existent (a se citi manualul din linux pentru comanda **vgextend**), insa acum un volume group nou si vom adauga un physical volume la el. Aceasta este procedura relativ usoara, o singura comanda, **vgcreate**, ce primeste numele noului volume group precum si physical volume ce va fi adaugat initial. De exemplu: **vgcreate vgdata /dev/vdb1**.

Este complet la alegerea utilizatorului ce nume alege pentru volume groups. In acest exemplu am ales sa incep numele cu vg, ceea ce le face usor de gasit daca vreodata avem foarte multe.

Dupa ce s-a creat volume group putem executa comanda **vgs** pentru a lista volume groups existente.

2.3.3 Pasi pentru a crea Logical Volumes

Acum ca s-a creat un volume group putem incepe sa cream logical volumes. Procedura este un pic mai complicata decat crearea physical volumes sau volume groups, deoarece sunt mai multe alegeri ce trebuie facute. Cand se creeaza un logical volumes trebuie specificate numele si marimea.

Marimea logical volume poate fi specificata prin valoare absoluta folosind optiunea **-L**. De exemplu folosind **-L 5G** se creeaza un logical volume cu marimea de 5GB. O alternativa este sa folosim optiunea **-l** pentru a specifica o marime relativa. De exemplu, se foloseste **-l 50% FREE** pentru a folosi jumatate din spatiul liber la nivelul volume group. Urmatorul argument este numele volume group la care va fi asignat logical volume nou creat si din care se va lua spatiu necesar. Un parametru optional este **-n** ce specifica numele noului logical volume. De exemplu, comanda **lvcreate -n lvvol1 -L 100M vgdata** creeaza un logical volume cu numele lvvol folosind 100MB spatiu din volume group cu numele vgdata (creat anterior).

2.3.4 Modificarea marimii unui Logical Volume

Motivul folosirii logical volumes in locul utilizarii directe a unor discuri este acela ca la nevoie putem foarte usor sa adaugam spatiu intr-un logical volume, in

timp ce daca de exemplu avem de adaugat in cazul discurilor trebuie sa adaugam un alt disc mai mare si sa copiem totul din discul initial, moment in care sistemul trebuie sa fie oprit. Daca de exemplu avem un folder /tenant1 in care tinem datele din baza de date pentru un client. Sa ne imaginal ca acesta trimite la un disc fizic de 100GB. Pentru a putea sa oferim mai mult spatiu in acest director va fi nevoie sa oprim baza de date, sa mutam toate datele pe un alt disc de 200GB de ex si apoi sa repornim sistemul. In cazul in care /tenant1 este un director la care a fost montat un logical volume putem folosi urmatoare procedura pentru a mari spatiul fara a opri sistemul\$

- Verificam ca in volume group din care face parte logical volume ce va fi extins se afla spatiu liber. Daca nu exista spatiu liber adauga un nou physical volume.
- Executam comanda **vgextend** pentru a extinde spatiul din logical volume.

Putem utiliza comanda **lvs** ce listeaza date despre logical volumes existente pentru a vedea spatiul liber al logical volume. Pentru detalii despre comanda **vgextend** a se executa comanda **man vgextend** ce afiseaza manualul din Linux.

2.4 Explicatii despre conventia de numire a LVM

Din acest moment putem incepe sa folosim logical volume nou creat. Pentru a face asta trebuie sa stim cum sa ne adresam acestuia. Entitatile LVM pot fi adresate in moduri multiple. Cea mai simpla metoda este sa adresam un logical volume precum: `/dev/vgname/lvname`. Astfel ca, daca am creat un logical volume cu numele `lvdata`, ce isi ia spatiu dintrun volume group cu numele `vgdata`, calea entitatii va fi `/dev/vgdata/lvdata`.

Device mapper poarta un rol important in denumirea LVM logical volumes. Device mapper (abreviat `dm`) este o interfata generica folosita de kernelul Linux pentru a adresa device de stocare. Device mapper este folosit de catre multiple sisteme: LVM volumes, dar si de software RAID etc. Aceste device sunt create doua locatii: device numerotate in secventa plasate in directorul `/dev`, precum `/dev/dm-0`, `/dev/dm-1` si asa mai departe. Deoarece aceste denumiri nu ofera nicio informatie despre deviceuri ci confuzeaza, se folosesc scurtaturi create in directorul `/dev/mapper`. Acestea folosesc sablonul `vgname-lvname` pentru nume. Astfel ca, daca avem device `/dev/vgdata/lvdata` putem sa ne adresam folosind si `/dev/mapper/vgdata-lvdata`. Cand lucram cum LVM logical volumes putem folosi oricare din aceste tipuri de adresare.

Database Point in time recovery Manager

Din punctul de vedere al unui dezvoltator software am observat in timpul experientei faptul ca pentru echipa efectiva de dezvoltare era destul de dificil de intretinut baza de date (ma refer chiar si la pornit/oprit, mai ales pentru juniori acest lucru putea pune probleme daca baza de date se afla instalata pe o platforma UNIX). In cazul in care se dorea aplicarea unei strategii de point in time recovery pentru restabilirea bazei de date la o stare existenta la un moment dat, era nevoie de interventia echipei de devOps, ceea ce se factura diferit la client, prin urmare trebuia convins in prealabil Product Owner-ul de necesitatea acelei interventii si planificat totul.

Astfel ca, am observat necesitatea existentei unei aplicatii care sa usureze operatiile asupra bazei de date, incepand cu pornirea/oprirea/verificarea statusului pana la intreaga operatie de backup si point in time recovery(intoarcere la o stare existenta la un anumit moment dat). Aceasta aplicatie trebuia sa fie intuitiva, usor de folosit chiar si pentru cineva care nu avea cunostinte de devOps, nefiind necesar accesul efectiv pe masina de UNIX unde se afla clusterul de baze de date. Este de la sine inteles ca aplicatia ar trebui sa detecteze toate bazele de date din clusterul de pe masina pe care este incarcata si ar trebui sa permita operatiile prezentate mai sus pentru fiecare dintre instanta in parte.

3.1 Arhitectura sistemului de fisiere

Am propus o solutie pentru bazele de date PostgreSQL folosite pe platforme UNIX (la momentul de fata). Pentru a prezenta solutia software am avut nevoie de un mediu UNIX pe care sa ruleze un numar de instante PostgreSQL. Am ales o masina CentOS, instanta EC2 oferita de Amazon. Sistemul de fisiere a fost organizat folosind LVM (logical volume manager) dupa cum urmeaza:

- am cerut 3 discuri de cate 50GiB, Amazon EBS;
- am creat cate un physical volume din fiecare din cele 3 discuri;
- am creat un volume group si le-am adaugat acestuia;
- am creat 6 logical volumes, 3 seturi de cate 2 logical volumes; un logical volume de 20GiB pentru datele din db si un logical volume de 10GiB pentru stocarea fisierelor WAL.

Au mai ramas 60GiB nealocati, acestia putand fi folositi ulterior. Daca se observa ca un client are nevoie de mai mult spatiu i se pot extinde logical volumes folosind din spatiul nealocat, acesta fiind mare avantaj al LVM. Bazele de date nu vor fi intrerupte cat timp se face extinderea logical volumes.

Am setat astfel sistemul de fisiere pentru 3 instante diferite de baze de date PostgreSQL denumite lv-pgdata-1, lv-pgdata-2 si lv-pgdata-3. In *Figura 1* este prezentata diagrama sistemului de fisiere.

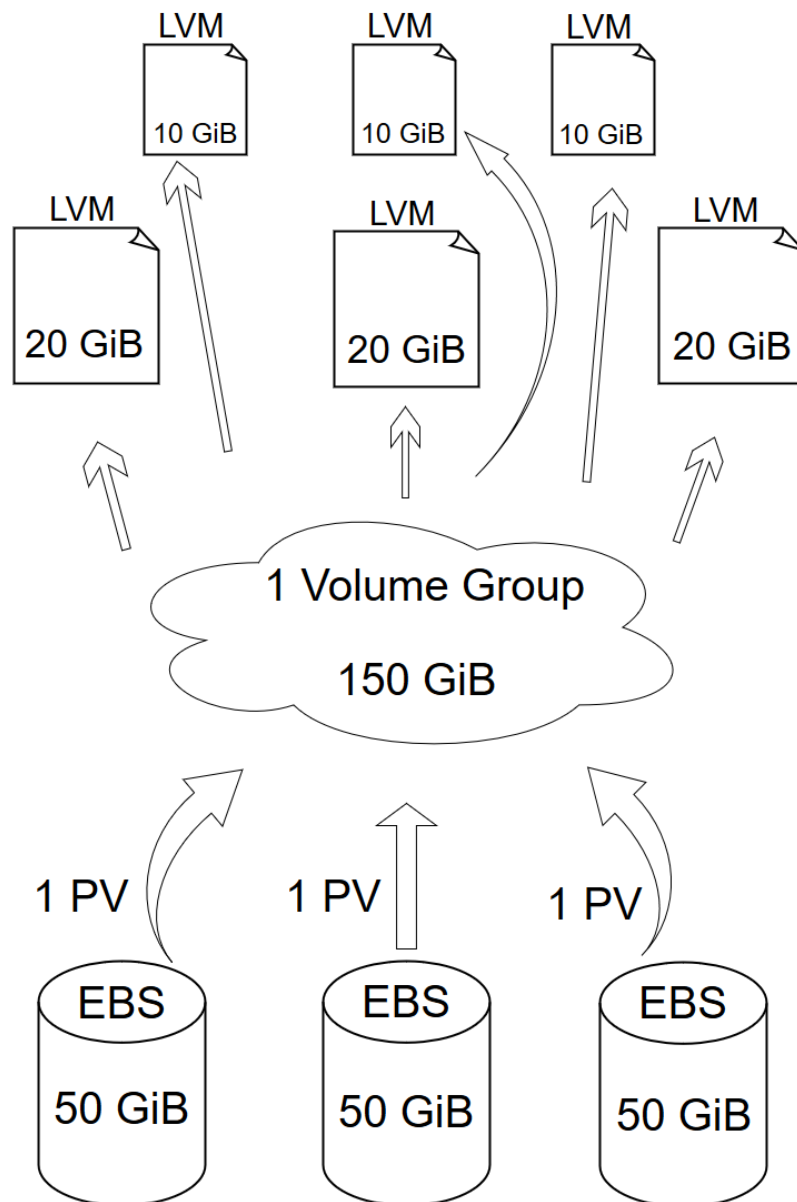


Figura 1: Diagrama sistemului de fisiere

DevOps as a Service

În cele ce urmează voi prezenta soluția software propusă pentru managementul bazelor de date în maniera descrisă anterior. Voi trece în revistă arhitectura software, urmând apoi să explic fiecare din componentele aplicației, precum și modul de interacționare a acestora.

4.1 Arhitectura Sistemului Software

Sistemul este împărțit în 3 componente importante:

- Scripturi UNIX ce interacționează cu instanțele de baze de date;
- RESTful API folosit pentru a expune funcționalitatea oferită de scripturi prin intermediul protocolului HTTP;
- Aplicație Web ce comunică prin intermediul protocolului HTTP cu API-ul RESTful expus.

În *Figura 2* este prezentată diagrama arhitecturii software.

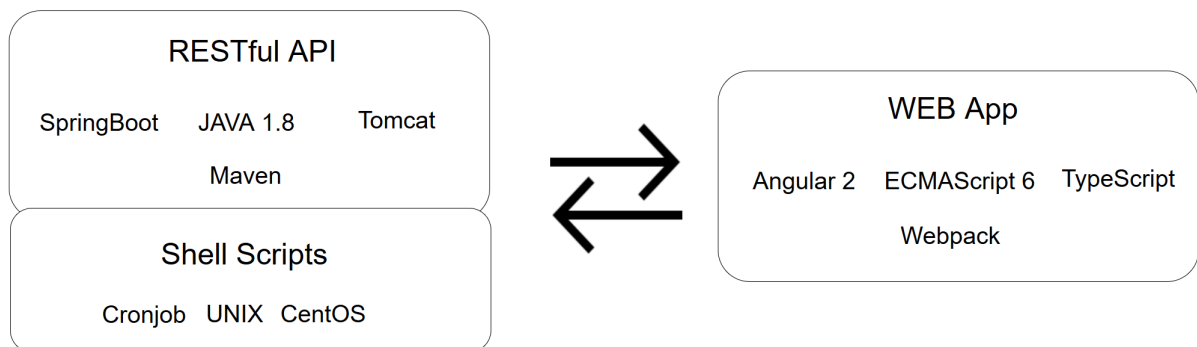


Figura 2: Diagrama arhitecturii software

4.2 Scripturile UNIX

Functionalitatea de baza a aplicatiei consta in interactionarea cu instantele de baze de date. Pentru asta am facut 3 scripturi ce au urmatoarele responsabilitati:

1. **pg-script.sh** pentru interactionarea cu instantele (pornire, oprire, interogare status);

```
1 #!/bin/bash
2
3 name="$(basename $0)"
4 pgdata="$(find /pgdata -type d -name $name)"
5 logFile="$pgdata/pg_log/startup.log"
6
7 if [ ! -f $logFile ]
8 then
9     "/bin/touch" "$logFile"
10 fi
11
12 "/usr/lib64/pgsql95/bin/pg_ctl" $1 "-D" $pgdata "-l" $logFile
```

Cod sursă 1: pg-script.sh

2. **basebackup.sh** pentru crearea de basebackup a unei instante (checkpoint de la care sa se poate incepe procesul de restaurare);

```
1 #!/bin/bash
2
3 if [ "$#" -ne 2 ]
4 then
5     echo "Illegal number of parameters. 2 params: ip port"
6     exit 1
7 fi
8
9 name="$(basename $0)"
10 pgdata="$(find /pgdata -type d -name $name)"
11 if [ ${#pgdata} == 0 ]
12 then
13     echo "tenant with name $name has no pgdata folder"
14     exit 2
15 fi
16
17 backup="$(find /backup -type d -name $name)"
18 if [ ${#backup} == 0 ]
19 then
20     echo "tenant with name $name has no backup folder"
21     exit 3
22 fi
23
24 basebackup="$backup/basebackup"
25 if [ -d $basebackup ]
26 then
27     rm -rf $basebackup/*
28     rm -rf $basebackup/.[a-z]*
29 else
30     mkdir $basebackup
31     chmod 700 $basebackup
32 fi
33
34 "pg_basebackup" "-h" $1 "-p" $2 "-D" $basebackup
```

Cod sursă 2: basebackup.sh

3. recovery.sh pentru initierea procesului de restaurare a bazei de date

```
1  #!/bin/bash
2
3  if [ "$#" -ne 2 ]
4  then
5      echo "Illegal number of parameters. 2 params: date e.g.
6          '2017-06-05 14:00:00 EEST' true/false if you want to start the db also"
7      exit 1
8  fi
9
10 name="$(basename $0)"
11 pgdata="$(find /pgdata -type d -name $name)"
12 if [ ${#pgdata} == 0 ]
13 then
14     echo "tenant with name $name has no pgdata folder"
15     exit 2
16 fi
17
18 backup="$(find /backup -type d -name $name)"
19 if [ ${#backup} == 0 ]
20 then
21     echo "tenant with name $name has no backup folder"
22     exit 3
23 fi
24
25 basebackup="$backup/basebackup"
26 if [ ${#basebackup} == 0 ]
27 then
28     echo "tenant with name $name has no basebackup folder"
29     exit 4
30 fi
31
32 wal="$backup/wal"
33 if [ ${#wal} == 0 ]
34 then
35     echo "tenant with name $name has no WAL files folder"
36     exit 5
37 fi
38
39 # TO DO: make sure that the server instance is stopped
40 "/home/ec2-user/disertatie/server/$name" "stop"
41
42 # TO DO: backup logs and other things of interes before empty pgdata dir
43 rm -rf $pgdata/*
44 rm -rf $pgdata/.[a-z]*
45
46 "cp" "-a" "-r" "$basebackup/." $pgdata
47
48 result="restore_command = 'cp $wal/%f %p'\n
49 recovery_target_time = '$1'"
50
51 echo -e $result > $pgdata/recovery.conf
52
53 if [ $2 == true ]
54 then
55     # Start the db server which will trigger recovery
56     "/home/ec2-user/disertatie/server/$name" "start"
57 fi
```

Cod sursă 3: recovery.sh

Toate cele trei scripturi au fost gandite cu scopul de a se crea symbolic links (scurtaturi) catre ele. Acestea vor avea numele instantei bazei de date, iar scriptul principal foloseste numele fisierului scurtatura pentru a accesa datele instantei respective. Scopul este acela de a crea un singur set de scripturi, iar daca avem n instante de baza de date doar creem symbolic links catre scripturi, fiecare cu numele instante de baza de date respective.

4.3 Server-ul: RESTful API

Pentru partea de server (backend) a fost aleasa solutia crearii unui RESTful API cu scopul de a se expune functionalitatea oferita de catre scripturile bash prezentate anterior. Am utilizat Java Enterprise Edition, framework-ul Spring, ce faciliteaza crearea unui API prin intermediul librariilor SpringWeb.

Un alt motiv pentru care am ales Spring a fost suportul oferit pentru a injecta dependinte urmand sablonul de proiectare Dependency Injection. Astfel, am decuplat controller-ul REST ce primeste cererea web de catre serviciul care o executa. Se poate injecta orice implementare de serviciu in controller, iar codul acestuia nu trebuie schimbat. Versiunea de Java utilizata este 1.8.

```
1 package pgapp.api;
2
3 import org.springframework.beans.factory.annotation.Autowired;
4 import org.springframework.web.bind.annotation.GetMapping;
5 import org.springframework.web.bind.annotation.PostMapping;
6 import org.springframework.web.bind.annotation.RequestMapping;
7 import org.springframework.web.bind.annotation.RequestParam;
8 import org.springframework.web.bind.annotation.RestController;
9 import pgapp.dto.DatabaseManagementDetailsObject;
10 import pgapp.service.DatabaseManagementService;
11
12 import java.util.List;
13
14 /**
15  * @author Stefan Stan on 06.06.2017.
16  */
17 @RestController
18 @RequestMapping("/pgapp")
19 public class DatabaseManagementAPI {
20
21     @Autowired
22     private DatabaseManagementService dbService;
23
24     @GetMapping("/server")
25     public List<String> getDBsNames() {
26         return dbService.getDBsNames();
27     }
28
29     @PostMapping("/server")
30     public DatabaseManagementDetailsObject changeDB(
31         @RequestParam("dbServerName") String dbServerName,
32         @RequestParam("action") String action
33     ) {
34         return dbService.changeDB(dbServerName, action);
35     }
36
37     @PostMapping("/basebackup")
38     public DatabaseManagementDetailsObject basebackupDB(
39         @RequestParam("dbServerName") String dbServerName
40     ) {
41         return dbService.basebackupDB(dbServerName);
42     }
43
44     @PostMapping("/pitr")
45     public DatabaseManagementDetailsObject recoverDB(
46         @RequestParam("dbServerName") String dbServerName,
47         @RequestParam("dateTime") String dateTime
48     ) {
49         return dbService.recoverDB(dbServerName, dateTime, true);
50     }
51 }
```

Cod sursă 4: DatabaseManagementAPI.java

Pentru managementul dependintelor si crearea artefactului s-a folosit Maven.

Se ofera ca si exemplu fisierul **pom.xml**:

```
1  <?xml version="1.0" encoding="UTF-8"?>
2  <project xmlns="http://maven.apache.org/POM/4.0.0"
3      xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
4      xsi:schemaLocation="http://maven.apache.org/POM/4.0.0
5          http://maven.apache.org/xsd/maven-4.0.0.xsd">
6      <modelVersion>4.0.0</modelVersion>
7
8      <parent>
9          <groupId>org.springframework.boot</groupId>
10         <artifactId>spring-boot-starter-parent</artifactId>
11         <version>1.5.2.RELEASE</version>
12     </parent>
13
14     <groupId>ro.stefanstan.apps</groupId>
15     <artifactId>pgApp</artifactId>
16     <version>1.0.0.1-SNAPSHOT</version>
17     <packaging>jar</packaging>
18
19     <properties>
20         <java.version>1.8</java.version>
21     </properties>
22
23     <dependencies>
24         <dependency>
25             <groupId>org.springframework.boot</groupId>
26             <artifactId>spring-boot-starter-web</artifactId>
27         </dependency>
28         <dependency>
29             <groupId>org.springframework.boot</groupId>
30             <artifactId>spring-boot-starter-test</artifactId>
31             <scope>test</scope>
32         </dependency>
33         <dependency>
34             <groupId>com.jayway.jsonpath</groupId>
35             <artifactId>json-path</artifactId>
36             <scope>test</scope>
37         </dependency>
38     </dependencies>
39
40     <build>
41         <plugins>
42             <plugin>
43                 <groupId>org.springframework.boot</groupId>
44                 <artifactId>spring-boot-maven-plugin</artifactId>
45             </plugin>
46         </plugins>
47     </build>
48
49     <repositories>
50         <repository>
51             <id>spring-releases</id>
52             <url>https://repo.spring.io/libs-release</url>
53         </repository>
54     </repositories>
55     <pluginRepositories>
56         <pluginRepository>
57             <id>spring-releases</id>
58             <url>https://repo.spring.io/libs-release</url>
59         </pluginRepository>
60     </pluginRepositories>
61 </project>
```

Cod sursă 5: pom.xml

4.4 Interfata grafica cu utilizatorul

Scopul proiectului este crearea unei aplicatii care sa faciliteze manipularea serverelor de baze de date PostgreSQL de catre dezvoltatori care nu trebuie sa aiba acces la masinile pe care se afla bazele de date/ cunostinte de UNIX. In acest sens, avand deja server-ul de Java care expune prin servicii Web functionalitatile respective, trebuie apelate si oferite intr-o maniera cat mai intuitiva.

Ca si framework am folosit Angular 2, limbajul de programare fiind Javascript, standardul ECMAScript6, impreuna cu TypeScript. Angular 2 are urmatoarele avantaje:

- este un framework open-source JavaScript construit și întreținut de Google, ceea ce usurează dezvoltarea aplicațiilor web.
- standardizează aplicațiile la nivel de client (browser), oferind o structură robustă și ușor de implementat în dezvoltarea site-urilor.
- este alegerea potrivită pentru orice aplicație web, în special dacă se preferă efectele vizuale spectaculoase. Rezultatul este un site fluid și rapid. Acesta este motivul pentru care se pretează SPA-urilor (Single Page Application) și aplicațiilor web destinate dispozitivelor mobile.
- se integrează ușor cu Bootstrap, permite crearea de aplicații responsive care pot fi accesate atât din browserul calculatorului cât și de pe dispozitivele mobile. Astfel, este ușor să creezi o singură aplicație web care să se vadă perfect pe orice dispozitiv de pe care este accesată – PC, tabletă, telefon, plasmă.

Ca si manager de pachete am utilizat **npm** de la NodeJS.

4.5 Setarea aplicatiei in Amazon Web Services

Avand ca si obiectiv crearea unui mediu cat mai apropiat de cel aflat pe masinile din productie, am inchiriat un server EC2 de la Amazon.

Sistemul de fisiere a fost stabilit conform descrierii din Capitolul 2. Am inchiriat 3 discuri EBS(Elastic Block Storage) pe care le-am configurat precum a fost prezentat, pentru a seta 3 instante de baza de date.

Odata ce sistemul de fisiere a fost configurat, am creat scripturile ce vor interactiona cu bazele de date. Urmatorul pas a fost pornirea celor doua servere (backend si frontend). In timpul acestui proces a fost necesara configurarea unui rol de securitate oferit masinii de calcul EC2, pentru a permite traficul pe portul 3001 utilizat de aplicatia frontend.

Ghid de utilizare a aplicației

Solutia propusa pentru a rezolva problema prezentata in capitolele precedente este o aplicatie web folosita pentru managementul bazelor de date.

Aceasta ofera functionalitatile de vizualizare a starii instantelor de baze de date prin apasarea butonului de "Refresh". Pentru a restaura starea bazei de date este necesara crearea unui checkpoint, iar pentru aceasta se apasa butonul "Create basebackup".

Odata existent un checkpoint, se poate restaura baza de date respectiva prin alegerea unei date/ore folosind componenta vizuala de tip calendar, urmand sa se apese butonul "Revert".

Directii de dezvoltare viitoare

Aplicatia propusa pentru a rezolva necesitatea restaurarii bazelor de date la anumite momente in timp poate fi folosita doar pentru baza de date PostgreSQL.

Astfel ca, o prima idee de dezvoltare viitoare ar fi investigarea suportului oferit alte tipuri de baze de date Oracle, MySql etc. urmand ca apoi sa creez un design abstract in ceea ce priveste baza de date si care sa fie integrat cu solutia existenta oferita pentru PostgreSQL.

Cat despre interactiunea cu utilizatorul, am folosit ultimele tehnologii actuale Angular 2, EcmaScript 6 etc. insa din punct de vedere al design, ar putea fi imbunatatit in ceea ce priveste html si css existent.

Concluzii

Am creat o aplicatie pornind de la necesitatea prezentata in capitolul *Introducere*. Scopul acesteia este de a oferi o metoda facila de a restaura o baza de date la un moment de timp oferit de catre utilizator. Acesta nu necesita cunostinte de baze de date, sau de devOps (nu e nevoie sa interactioneze cu scripturi de UNIX).

Pentru a realiza cele mentionate anterior am studiat modalitatea bazei de date PostgreSQL de management intern a transactiilor si am utilizat-o intro maniera ce imi permite restaurarea bazei de date la un moment in timp cu o granularitate fina, la secunda.

De asemenea, am studiat modalitatea sistemului de operare UNIX ce permite creare de partitii virtuale ce pot fi oricand modificate in ceea ce priveste marimea fara a inchide baza de date. Am combinat aceasta functionalitate a sistemului de operare impreuna cu solutia propusa, pentru a realiza cele prezentate cu un timp minim de inaccesibilitate a bazei de date. *Capitolul 3* prezinta in detaliu aceste aspecte.

Am propus o arhitectura a sistemului de fisiere bazat pe studiul prezentat anterior si am utilizat-o in implementarea propriu-zisa. Cat despre aceasta, am folosit serviciile Amazon Cloud (EC2, EBS) pentru a face deploy solutiei (scripturi ce interactioneaza cu baza de date, aplicatie backend si aplicatie frontend).

In opinia mea, prin studiul realizat si prin intermediul solutiei propuse am reusit sa realizez ce mi-am propus si am oferit o modalitate de management a bazelor de date si de restaurare a acestora la starea existenta un moment de timp. Prin urmare am realizat o aplicatie usor de utilizat de catre persoane ce nu au cunostinta de baze de date, sistem de operare UNIX.

Bibliografie

- [1] Sander van Vugt, Red Hat® RHCSA™/RHCE® 7 Cert Guide: Red Hat Enterprise Linux 7 (EX200 and EX300)
- [2] E. Gamma, R. Helm, R. Johnson, and J. Vlissides. *Design Patterns: Elements of Reusable Object-Oriented Software*. Addison-Wesley, 2005.
- [3] Postgres Documentation
<https://www.postgresql.org/docs/9.5/static/continuous-archiving.html>
- [4] Craig Walls, Spring in Action, Third Edition
- [5] Martin Fowler. *GUI Architectures*. 18 July 2006.