

Mailtrace – E-Mail-Recherche für Helpdesks

Heinlein – Professional Linux Support GmbH

Stefan Neben

s.neben@heinlein-support.de
030 / 40 50 51 - 0

E-Mail-Recherche, was ist damit gemeint?

- ▶ Kurzgefasst das Auswerten der Logdateien von Mailservern
- ▶ Situation bei ISPs, wenn manuell ausgewertet werden muss
 - ▶ Kunde ruft an und erkundigt sich zum Verbleib einer E-Mail
 - ▶ Einloggen des Supporters auf den Mailserver bzw. den zentralen Logserver
 - ▶ (Bzw. Beauftragung des verantwortlichen Administrators)
 - ▶ `grep` über die Logfiles, bis gewünschte Informationen gefunden sind
 - ▶ Weitere `grep`-Suchen bei komplexeren Mailsetups
- ▶ **FAZIT:** So oder so ist diese Aufgabe zeitraubend!

Wie kann man diesem Problem entgegentreten?

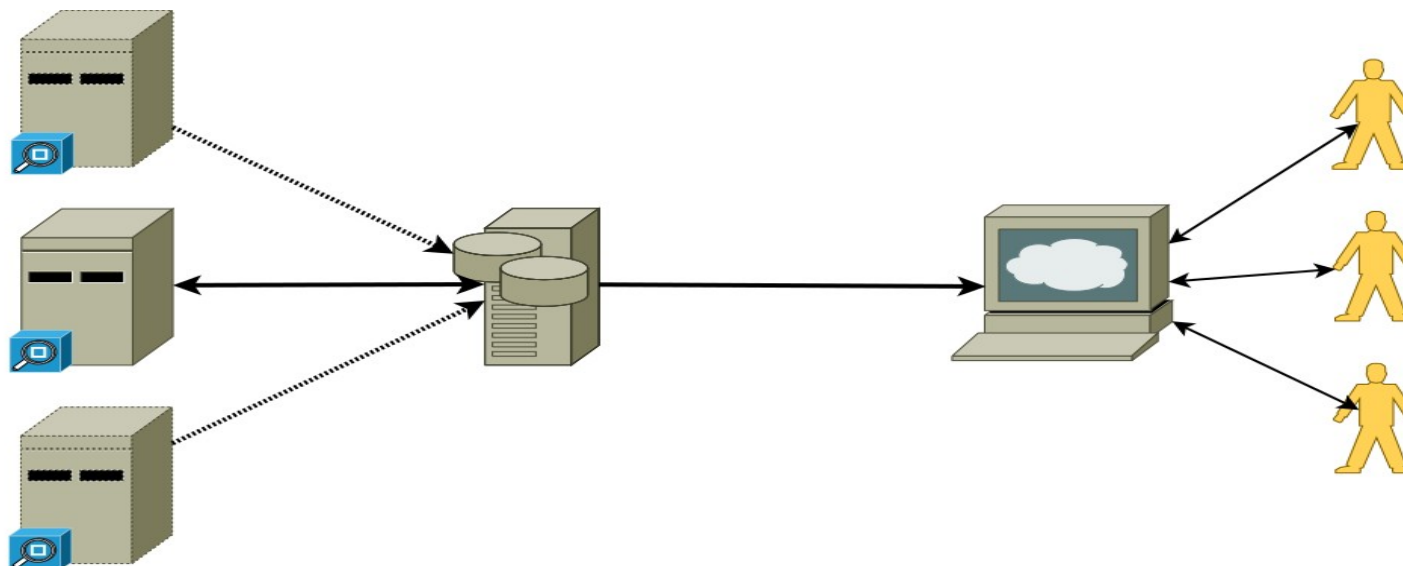
- ▶ Notwendige/interessante Daten müssen in eine Datenbank
- ▶ **Warum?**
 - ▶ Daten liegen dann strukturiert/gefiltert vor
 - ▶ Durch das Datenbank-Schema existiert eine vereinheitlichte Schnittstelle
 - ▶ Somit kann einfacher eine eigene Programmlogik zur Realisierung einer grafischen Oberfläche verwendet werden
 - ▶ Nutzerverwaltung ist dadurch einfacher (Zugriff von Nicht-Administratoren)

Probleme bei der Problemlösung

- ▶ Die Datenbank ist der Flaschenhals
 - ▶ Unter anderem geschuldet des Log-Aufbaus von Postfix
 - ▶ Je größer die bereitgestellte Datenmenge, desto langsamer die Datenbank
- ▶ Der Datenbestand soll lückenlos sein!
 - ▶ Was passiert bei Ausfall der Datenbank
 - ▶ Abfangen von Fremdeinwirkungen (Logrotate, etc.)
- ▶ Duplikate sollen weitestgehend vermieden werden

Grundsätzliche Herangehensweise zur automatisierten Mail-Log-Auswertung

- ▶ Drei Hauptkomponenten
 - ▶ Ein Mail-Logparser als Daemon auf dem/den Mailserver(n)
 - ▶ Eine zentrale Datenbank
 - ▶ Eine WebGUI als Interface für den Endanwender



Der Mailtrace-Daemon

- ▶ Geschrieben in Perl
- ▶ Verwendet hauptsächlich Standard-Module, welche über das Paketmanagement angeboten werden (Debian, SuSE)
- ▶ Durchläuft beim Start die Logdatei komplett und schaltet am Datei-Ende in einen „Tail-Modus“
- ▶ Ist der einzige schreibende Part für die Datenbank, auch für Löschvorgänge

Der Mailtrace-Daemon: Problem Nr. 1

- ▶ Loggingverhalten von Postfix schlecht für die Gesamtperformance

- ▶ Schreibt für einen Vorgang mehrere Logzeilen

- ▶ Ein Beispiel:

```
postfix/smtpd[22503]: connect from host[10.20.30.40]
postfix/cleanup[22504]: 4FFAD12389BC: message-id=<20110306110556>
postfix/qmgr[18150]: 4FFAD12389BC: from=<mail@mail.de>, size=6815, nrcpt=1 (queue active)
postfix/pipe[22512]: 4FFAD12389BC: to=<name@mail.de>, relay=dovecot, delay=0.25,
    delays=0.03/0.02/0/0.2, dsn=2.0.0, status=sent (delivered via dovecot service)
```

- ▶ Interessante Daten über vier Zeilen verteilt

- ▶ Der erste Eintrag wird mit unvollständigen Daten mittels INSERT in DB geschrieben
 - ▶ Jede weitere Zeile sorgt mittels UPDATE zur weiteren Vervollständigung
 - ▶ **Wie kann nun dieser Overhead minimiert werden?**

Der Mailtrace-Daemon: Problem Nr. 1

- ▶ **Der Memcached schafft Abhilfe!**
 - ▶ Caching-Daemon die „lokale“ Maschine oder über Netzwerk
 - ▶ Sehr performant
 - ▶ Verhalten lässt sich gut konfigurieren (u.A. Vorhaltezeit und Größe des Caches)
 - ▶ Bibliotheken für fast alle Programmiersprachen verfügbar
- ▶ Durch die Verwendung des Memcached lassen sich die Datenbank-Zugriffe auf mindestens 2 minimieren
- ▶ **Warum mindestens 2?**
 - ▶ Geschuldet dem derzeit verwendeten Datenbankschema (dazu später mehr)
 - ▶ Aber auch dem geschuldet, wie E-Mails beim Versand verarbeitet werden

Der Mailtrace-Daemon: Problem Nr. 2

- ▶ Was wenn die Datenbank nicht erreichbar ist?
 - ▶ Wie große Lücken vermeiden?
- ▶ Daemon muss daher wie folgt ausgelegt sein:
 - ▶ Genaues Pausieren an der aktuellen Abarbeitungsstelle bei Verbindungsverlust
 - ▶ In eine Wiederverbindungs-Schleife gehen, um die Datenbank-Verbindung wieder herzustellen
 - ▶ Wiederaufnahme des Parsing-Vorgangs nach wieder hergestellter Verbindung
- ▶ In einem solchen Fall nimmt die Wahrscheinlichkeit von Fremdeinwirkung zu (z.B. logrotate)
 - ▶ **Gefahr von großen Lücken entsteht!**

Der Mailtrace-Daemon: Problem Nr. 3

- ▶ Logrotate rotiert Datei weg, obwohl der Daemon noch nicht das Ende erreicht hat
 - ▶ Je nach Umstand erhöht sich die Gefahr von großen Lücken
- ▶ Die Lösung ist ein Problem
 - ▶ Einige Administratoren werden das Problem sicherlich kennen, welches in diesem Fall die Lösung darstellt
 - ▶ Prozess friert z.B. ein, hat eine Datei geöffnet und diese wird gelöscht
 - ▶ Doch es wird kein Plattenplatz freigegeben, warum?
 - ▶ Verzeichnis-Zuordnung wird entfernt und die Datei ist nicht mehr „sichtbar“
 - ▶ Der wirklich belegte Platz wird aber erst freigegeben, wenn kein Zugriff mehr
- ▶ **Diese Besonderheit hilft!**

Der Mailtrace-Daemon: Problem Nr. 3

- ▶ Daemon muss daher multi-instanzfähig sein
 - ▶ Verwendung von eindeutigen PID-Dateien (`mailtraced.<timestamp>.pid`)
- ▶ Bei `logrotate` wird für das `postrotate` folgendes konfiguriert
 - ▶ Senden von `SIGTERM` an alle derzeit laufenden Instanzen
 - ▶ Starten einer neuen Instanz
 - ▶ Der Daemon fängt das `SIGTERM` ab, beendet sich aber nicht sofort
 - ▶ Die rotierte alte Logdatei ist zwar nicht mehr sichtbar, aber noch vorhanden:

```
# lsof | grep 29095
...
mailtrace 29095 root 3r REG 8,2 431698710 312189 /var/log/mail-20110310 (deleted)
...
```
 - ▶ Erst wenn das dortige Ende erreicht ist, soll sich der Daemon beenden

Der Mailtrace-Daemon: Problem Nr. 4

- ▶ Daemon muss aus bestimmten Gründen neu gestartet werden
 - ▶ Duplikate sollen hierbei vermieden werden
- ▶ Wie wird verhindert, dass ein bereits geparster Eintrag erneut in die Datenbank geschrieben wird?
 - ▶ **Möglichkeit 1:**
 - ▶ Unique-Key über mehrere Spalten in der Datenbank
 - ▶ Nachteil ist ein erhöhter Platzverbrauch auf dem Datenbank-Server
 - ▶ **Möglichkeit 2:**
 - ▶ Man nimmt eine erhöhte Anzahl von `SELECT`-Anfragen zur Prüfung in Kauf
 - ▶ Zumindest, bis das Dateiende erreicht ist
 - ▶ Wenn man kleine bis große Setups bedienen möchte, der beste Kompromiss

Der Mailtrace-Daemon: Problem Nr. 4

- ▶ Umsetzung der Dopplungsprüfung mittels `SELECT`-Anfragen
 - ▶ Aufteilung des Betriebs-Modi des Daemon in zwei Bereiche
 - ▶ **Modus 1:**
 - ▶ `SELECT` auf ausgewählte Attribute, bevor der ermittelte Datensatz eingefügt wird
 - ▶ Werte die sonst den Multi-Column-Unique-Key darstellen würden
(*Mail-Adresse, Queue-ID, Antwort des Servers, Server, Zeit*)
 - ▶ **Modus 2:**
 - ▶ Direktes Schreiben des Datensatzes ohne weitere Überprüfung
- ▶ Umsetzung im Daemon durch objektorientierten Ansatz

Der Mailtrace-Daemon: Problem Nr. 5

- ▶ Ist an sich **nicht** ein Problem des Daemons
 - ▶ Sollte aber über ihn gesteuert werden, da er die steuernde Instanz ist
- ▶ Damit der Datenbank-Server immer genügend Speicherplatz hat, müssen alte Daten entfernt werden
- ▶ Die Funktion `–max-age=<days_back>` wird auf **einem** ausgewählten Server in die Crontab eingetragen
 - ▶ Mit dieser Funktion werden alle Einträge älter als die angegebenen letzten Tage aus der Datenbank entfernt

Der Mailtrace-Daemon: Weitere Tuning-Möglichkeiten

- ▶ Der Großteil der Arbeit wird über Regex erledigt
 - ▶ Daher Pre-Compiling aller genutzten regulären Ausdrücke
- ▶ Handhabung der Daten
 - ▶ Häufige Verwendung von Referenzen (Zeigern), damit nur Speicheradressen zwischen den Funktionen ausgetauscht werden müssen
- ▶ Eigenes Cache-Handling der Datenbank-Handles
 - ▶ Durch den objektorientierten Ansatz können Daten zur Laufzeit besser vorgehalten werden
 - ▶ Somit müssen nur noch die Daten übertragen werden

Der Datenbank-Server

- ▶ Speicherplatz sollte je nach Setup großzügig ausgelegt werden
- ▶ Hauptverwendung von PostgreSQL, aber auch MySQL ist möglich
- ▶ Für eine bessere Performance:
 - ▶ Sollte der Datenbank-Server dediziert sein
 - ▶ Optimal daher ein eigener Server für jeden Daemon (logisch), für den Datenbank-Server und für den Webserver
 - ▶ Nur in vorher gut durchdachten kleineren Setups, können Teile auf einen Server zusammengelegt werden

Das Datenbank-Schema

- ▶ Was gilt es bei der Verarbeitung von E-Mails rein technisch zu beachten?
 - ▶ Es gibt immer einen Verbindungsaufbau und einen Absender
 - ▶ Aber es kann mehrere Empfänger geben
 - ▶ Der Mailserver behandelt jede E-Mail für jeden Empfänger einzeln
- ▶ Deshalb Aufteilung des Schemas in `connectlog`- und `deliverlog`-Sektionen
 - ▶ Ein `INSERT` für jeden Connect, also Absender
 - ▶ Ein `INSERT` für jedes Deliver, also Empfänger
 - ▶ Im Minimal-Fall (ein Absender und ein Empfänger) daher ein Zugriff auf die Datenbank mit mindestens 2 `INSERT`-Statements

Das Datenbank-Schema

- ▶ Damit die spätere Suche durch das Frontend effektiv ist, muss aber ein kleiner Bruch in der Logik hingenommen werden
 - ▶ Deliver-Einträge werden in die entsprechende Tabelle mit dem Zeitstempel des Connects geschrieben
- ▶ Doch die Tabellen können mit der Zeit eine enorme Größe erlangen und die Performance kann leiden
 - ▶ Ein Index hilft zwar, aber das reicht noch nicht!
 - ▶ Eine Aufteilung der Tabellen in Tagen schafft Abhilfe
 - ▶ Viele Suchanfragen konzentrieren sich hauptsächlich auf die letzten 2-3 Tage
 - ▶ *Deshalb muss oft nicht bis in die Tiefe gesucht werden*
 - ▶ Kleinhalten der Datenmenge notwendig (immer alte Einträge löschen)

Das Nutzerinterface

- ▶ Was aber kostet den Helpdesk die meiste Zeit?
 - ▶ Der Anruf des Kunden selber
- ▶ Die häufigsten Antworten:
 - ▶ Greylisting (bei sehr Ungeduldigen)
 - ▶ Temporäre Fehler der Gegenstelle
 - ▶ Hart geblockte E-Mails (Bounce-Meldung ist im Spam-Verdachtsordner gelandet)
- ▶ Wie kann man den Kunden dazu bringen, erst gar nicht nach dem Verbleib einer E-Mail zu fragen?
 - ▶ Wenn er die Möglichkeit zur eigenen Recherche hat
 - ▶ Daher ist für die Akzeptanz auch eine hohe Suchgeschwindigkeit enorm wichtig

Das Nutzerinterface

- ▶ Folgende Suchmöglichkeiten stehen im Frontend zur Verfügung

Absender (from:)	<input type="text"/>	Zeitraum von (Connect) ✕	<input type="text"/>
Empfänger (to:)	<input type="text"/>	Zeitraum bis (Connect) ✕	<input type="text"/>
Message-ID [?]	<input type="text"/>	Zeitraum von (Deliver) ✕	<input type="text"/>
Versandstatus [?]	--jeder-- ▼	Zeitraum bis (Deliver) ✕	<input type="text"/>
↑ Erweiterte Suche			
SMTP-Logmeldung [?]	<input type="text"/>	Bekommen von Server [?]	<input type="text"/>
Queue-ID (Connect) [?]	<input type="text"/>	Gesendet an Server [?]	<input type="text"/>
Queue-ID (Deliver) [?]	<input type="text"/>	Geloggt von Server [?]	--jeder-- ▼
➔ Neue Suche starten			<input type="button" value="Suchen"/>

Das Nutzerinterface

- ▶ Es kann nach allen gespeicherten Informationen gesucht werden (so sollte es ja auch sein, zumindest für den Admin)
- ▶ Erfahrungsgemäß sind aber für den Helpdesk oder den Nutzer selber „nur“ folgende Werte interessant:
 - ▶ Von wem kam die E-Mail
 - ▶ Zu wem wurde die E-Mail versandt
 - ▶ Wann ungefähr wurde die Mail verschickt (Zeit der Einlieferung)
- ▶ Durch die Partitionierung der Tabellen nach Tagen ist ein Suchvorgang auch bei großen Datenmengen kurz

Das Nutzerinterface

► Ein Ausschnitt einer Suchabfrage

25 ▼	pro Seite	Zeit/C ↑↓	Absender ↑↓	Empfänger ↑↓	Status ↑↓	1
-> 2011-02-21 06:22:33			v.test@heinlein-support.de	→ user@yahoo.de		
<- 2011-02-21 06:22:33			Message-ID: 1102220114573Y.12580@mail.heinlein-support.de		● Nachricht erfolgreich versendet ↓	Technische Details
-> 2011-02-28 11:21:13			spammer@spam.org	→ nutzer@heinlein-support.de		
<- 2011-02-28 11:21:13			Message-ID: 5298144515@spam.org		● rejected ↓	Technische Details
-> 2011-02-21 20:52:24			spam@spammer.de	→ info@heinlein-support.de		
<- 2011-02-21 20:52:24			Message-ID: unknown		● blocked ↓	Technische Details
-> 2011-02-21 23:54:02			info@heinlein-support.de	→ user@yahoo.de		
<- 2011-02-21 23:54:02			Message-ID: 1102220114573Y.12580@mail.heinlein-support.de		● Nachricht noch nicht versendet ↓	Technische Details
-> 2011-02-21 23:58:02			info@heinlein-support.de	→ user@yahoo.de		
<- 2011-02-21 23:58:02			Message-ID: 1102220114573Y.12580@mail.heinlein-support.de		● Nachricht erfolgreich versendet ↓	Technische Details
-> 2011-03-10 17:20:20			user@heinlein-support.de	→ chef@example.com		
<- 2011-03-10 17:20:21			Message-ID: 387ec12815574aa2577055bfc0ba0b1b@127.0.0.1		● Nachricht erfolgreich versendet ↓	Technische Details

Das Nutzerinterface

- ▶ In jedem einzelnen Suchergebnis stehen folgende Möglichkeiten zur Verfügung
 - ▶ Zusammenfassung aller Mails mit bestimmter Message-ID
 - ▶ Aufruf einer Übersichtsseite, welche alle gespeicherten Informationen zur jeweiligen E-Mail ausgibt
 - ▶ Eine Aufklapp-Funktion mit:
 - ▶ Mailfluss-Zusammenfassung (noch nicht vollständig implementiert)
 - ▶ Eine „menschenslesbare“ Meldung, was mit der E-Mail schlussendlich passiert ist

- ▶ Weitere Details in der Vorführung ...

Eine kleine Demonstration

- ▶ Demonstration des Daemons
- ▶ Demonstration der Oberfläche

Heinlein Support hilft auch bei allen anderen Fragen rund um Linux:

▶ AKADEMIE

- ▶ Von Profis für Profis: Wir vermitteln die oberen 10% Wissen. Geballtes Wissen und umfangreiche Praxiserfahrung aus erster Hand.

▶ SUPPORT

- ▶ Wir sind das Backup für Ihre Linux-Administration: LPIC-2-Profis lösen im Heinlein CompetenceCall Notfälle, auf Wunsch auch in SLAs mit 24/7-Verfügbarkeiten.

▶ HOSTING

- ▶ Wenn Hosting kein Massengeschäft sein darf: Individuelles Business-Hosting mit perfekter Maintenance durch unsere Linux-Profis. Sicherheit und Verfügbarkeit werden bei uns groß geschrieben.