

---

# Eine Einführung in sed

Thomas Pircher, [http://www.tty1.net/sed-tutorium\\_de.html](http://www.tty1.net/sed-tutorium_de.html) <tehpeh@gmx.net>

Dieses Dokument steht unter der Creative Commons Attribution-Share Alike 3.0 Unported Lizenz [<http://creativecommons.org/licenses/by-sa/3.0/>].

v1.4, August 2008

## Zusammenfassung

Seit ich mir ein bisschen Zeit genommen habe mich in **sed** einzuarbeiten, hat dieses Tool mir sehr viele Aufgaben erleichtert. Es ist unglaublich, wie flexibel dieses Programm einsetzbar ist und welche komplexen Regeln man mit ein paar Zeichenkombinationen aufstellen kann, die für den Laien nur wie ein beliebiges Gekrösel von Hieroglyphen ausschauen. Außerdem macht es einfach Spaß, eigene Scripte zu erstellen oder die von anderen Leuten zu verstehen. Mit dieser Einführung will ich Lust auf **sed** machen und ich hoffe dass die Einarbeitung in das Programm nicht zur Frust ausartet.

## Inhaltsverzeichnis

Einführung .....	2
Über dieses Tutorium .....	2
Was ist sed .....	2
Das Programm .....	2
Dokumentationen .....	3
Grundlagen .....	3
Reguläre Ausdrücke .....	3
Regular Examples .....	4
Grenzen von REs .....	5
<b>sed</b> und REs .....	5
Erste Schritte mit sed .....	5
Ein- und Ausgabe .....	5
Kommandos .....	6
Adressen .....	6
Mehr Kommandos .....	7
Ein paar interessantere Beispiele .....	9
Probleme mit REs .....	9
Selektives Ersetzen .....	10
Mehrere Kommandos .....	11
<i>spaceballs</i> .....	11
Ergänzungen zum <i>pattern space</i> .....	11
Einmal <i>hold space</i> und zurück .....	12
Sprünge ( <i>branches</i> ) .....	13
Sprungkommandos .....	13
Andere Sprünge .....	14
Vermischtes .....	14
Dateien .....	14
Noch mehr Kommandos .....	14
To <b>sed</b> or not to <b>sed</b> ? .....	15
Andere Programme mit <b>sed</b> -Kommandos .....	15
Ein paar Beispiele .....	15
Entfernen von Kommentaren .....	16
<i>elleff</i> -Rücktransformation .....	18

Verschachtelte Klammern .....	18
Kurzreferenz .....	18
Adressen .....	18
Kommandos .....	19
Versionsgeschichte .....	20

# Einführung

## Über dieses Tutorium

Ich bin kein **sed**-Guru, deshalb könnte manches Problem sicher mit weniger Kommandos gelöst werden. Das ist aber nicht der Sinn dieser Einführung die versucht die Beispiele so nachvollziehbar wie nur möglich zu halten. Ich bin aber dankbar für jedes Feedback, für Anregungen, Kritik, Fehlermeldungen, Verbesserungsvorschläge und Blankoschecks.

Die hier vorgestellten Scripte wurden in der Regel getestet, doch kann ich nicht ausschließen, dass sich nicht doch hier und da einige Fehler eingeschlichen haben. In diesem Fall bitte nicht meckern sondern melden. Für Fragen stehe ich gerne zur Verfügung. Einfach eine eMail an die auf der ersten Seite angegebene Adresse schicken und in der Betreff-Zeile '*sed-tutorium*' unterbringen.

Neue Versionen dieses Dokumentes können von [www.tty1.net/sed-tutorium\\_de.html](http://www.tty1.net/sed-tutorium_de.html) [[http://www.tty1.net/sed-tutorium\\_de.html](http://www.tty1.net/sed-tutorium_de.html)] bezogen werden.

## Was ist sed

**sed** wurde 1973 oder 1974 von Lee E. McMahon geschrieben (mehr zur Geschichte der frühen UNIX-Kommandos auf <http://www.columbia.edu/~rh120/ch106.x09>). Der Name **sed** steht für Stream Editor, was andeuten soll, dass das Programm kein Editor im gewöhnlichen Sinn ist, sondern Zeile für Zeile von *stdin* (standard input, meistens die Konsole) in seinen *pattern buffer* liest, diesen nach bestimmten Regeln bearbeitet, und anschließend den bearbeiteten *pattern buffer* auf *stdout* (standard output, meistens der Bildschirm), ausgibt. Die Quelldatei bleibt dabei unverändert.

Zur Verdeutlichung der Arbeitsweise von **sed** hier ein kleines Beispiel:

```
sed -e '/root/d'
```

Von auf der Tastatur eingegebenen Zeilen werden jene weggefiltert, die den String 'root' enthalten. Das wurde so im Script '/root/d' angegeben. Zum Beenden des Programms drückt man **^D** (CTRL und **D**) was unter UNIX für Dateiende steht.

Trotz seines Alters wird **sed** immer noch verwendet, da es ein weites Spektrum von einfachen bis sehr komplexen Aufgaben erledigen kann. Die Stärke von UNIX liegt zum Teil auch darin, dass es über kleine hochspezialisierte Tools verfügt die, miteinander kombiniert, fast jede Aufgabe lösen können.

## Das Programm

Bedingt durch die weite Verbreitung von **sed** gibt es eine Reihe von Implementationen, die sich in einigen Details unterscheiden. Dieses Tutorium versucht, sich so weit als möglich an die standardisierte Version von **sed** zu halten. Dem Autor liegt GNU **sed** vor, und alle Beispiele wurden damit getestet. Im Text wird aber jeweils darauf hingewiesen, wenn Erweiterungen verwendet werden. Auf Pements sedfaq [<http://sed.sourceforge.net/sedfaq.html>] findet sich eine erschöpfende Liste der verschiedenen **sed**-Implementationen und deren Unterschiede.

Eine gute Hersteller-unabhängige Dokumentation stellt die Open Group [<http://www.opengroup.org/onlinepubs/007908799/xcu/sed.html>] bereit.

## Dokumentationen

Zu nennen ist natürlich die man-page zu **sed**, die aber eher für schon sattelfeste Benutzer verdaulich ist. Ausführlicher und systematischer ist die *info* page zu **sed**. Damit sollte man nach der Lektüre dieser Einführung keine Probleme mehr haben.

Eine Sammlung interessanter Fragen rund um **sed**, inklusive vieler Scripte sowie eine nicht enden wollende Liste von weiterführenden Links findet man auf der sedfaq [<http://sed.sourceforge.net/sedfaq.html>] von Eric Pement. Das Tutorial von Carlos Jorge Duarte (do it with sed [[http://www.mds.mdh.se/~dat95abs/sed\\_tutorial.txt](http://www.mds.mdh.se/~dat95abs/sed_tutorial.txt)]) ist sehr lesenswert, besonders wegen der vielen sehr gut dokumentierten und zum Teil trickreichen Scripte. Wer keine Dokumentation zu **sed** hat, findet in dem File eine kurze aber durchaus brauchbare Referenz. Sehr gut ist auch u-sedit [<ftp://ftp.cs.umu.se/pub/pc/u-sedit2.zip>] von Mike Arst, welches ein Tutorial und jede Menge Beispiele beinhaltet.

Spezifische Fragen (sofern nicht in den FAQs behandelt) kann man auch in den *dafür zuständigen* Newsgroups stellen. Man bekommt dort meistens eine fundierte und ausführliche Antwort - eine treffende Beschreibung des Problems und die Befolgung der jeweiligen Netiquette seitens des Fragenden vorausgesetzt. Eine zuständige deutschsprachige Newsgruppe könnte beispielsweise `de.comp.os.unix.shell` sein.

## Grundlagen

### Reguläre Ausdrücke

Reguläre Ausdrücke (RE, Regular Expressions) wurden von 1956 S.C. Kleene eingeführt und sie erwiesen sich als sehr effektiv, um Zeichenketten zu beschreiben.

REs werden dann verwendet, wenn man die *Form* einer Zeichenkette (String) angeben will - sie beschreiben also *Klassen* von Strings. Es ist zum Beispiel einfacher die Form der *natürlichen Zahlen* als "eine Zeichenkette, bestehend aus einer oder mehreren Ziffern aus der Menge {0,1,2,3,4,5,6,7,8,9}" zu definieren, als alle Zahlen von 0 bis unendlich aufzuzählen. Reguläre Ausdrücke sind eine Schreibweise, mit der man solche Klassen eindeutig beschreiben kann.

Man sagt eine RE passt auf eine Zeichenkette, wenn diese in der von der RE umrissenen Klasse enthalten ist. Häufig werden REs verwendet, um aus einem String einen Teilstring heraus zu picken, welcher von der RE beschrieben werden kann. Dabei gilt das Prinzip der längsten Übereinstimmung (*longest match*), was heißen soll dass dies der längste Teilstring ist, auf den die RE passt. In diesem Zusammenhang spricht man auch davon, REs sind *greedy*.

**Tabelle 1. Erweiterte Reguläre Ausdrücke (Extended Regular Expressions)**

Regulärer Ausdruck	Erklärung
x	das Zeichen 'x'
\x	wenn x in [abfnrtv] dann die ANSI C-Interpretation davon, sonst x selber, z.B. \<*
\123	das Zeichen mit oktalem ASCII-Code 123
\xe5	das Zeichen mit hexadezimalen ASCII-Code e5
.	jedes beliebige Zeichen außer \n (newline)
[xyz]	eine "character class": x ODER y ODER z
[ako-sP]	eine "character class" mit einer Bereichsangabe, also a ODER k ODER ein Character aus dem Bereich o BIS s ODER P
[^x-z]	eine "negated character class": Jedes Zeichen außer x bis z
(r)	die RE r selber
rs	die RE r gefolgt von der RE s
r s	die RE r ODER die RE s
r*	die RE r null oder mehrere male
r+	die RE r ein oder mehrere male
r?	die RE r null oder ein mal
r{2,6}	die RE r zwei bis sechs mal
r{2,}	die RE r zwei oder mehrere male
r{,6}	die RE r null bis sechs mal
r{4}	die RE r genau vier mal
^r	die RE r am Anfang der Zeile
r\$	die RE r am Ende der Zeile
[:str:]	mit str eine der folgenden Bezeichner: <i>alnum</i> , <i>alpha</i> , <i>blank</i> , <i>cntrl</i> , <i>digit</i> , <i>graph</i> , <i>lower</i> , <i>print</i> , <i>punct</i> , <i>space</i> , <i>upper</i> , <i>xdigit</i> dann die betreffende Charakterklasse. Siehe <i>ctype(3)</i> für Details.

Für detailliertere Informationen siehe die man-page *regex(7)* oder, falls nicht vorhanden die man-page zu *awk(1)*, welche lange Zeit auch als die Referenz für REs galt, oder *flex(1)* oder die Onlinedokumentation der Open Group [<http://www.opengroup.org/onlinepubs/007908799/xbd/re.html>].

## Regular Examples

(x|y|z) ist äquivalent zu [xyz] und (a|b) ist äquivalent zu (b|a). Die RE (B|Dr)a1{2} passt sowohl auf "Ball" als auch auf "Drall". Die RE ^#. \* passt auf alle Zeilen, die mit einem '#' anfangen.

Ist die RE ^# äquivalent zur vorhergehenden? Zur reinen Mustersuche passen beide REs auf die selben Zeilen, der Unterschied kommt dann zu Tage, wenn man die gefundenen Muster weiterverarbeiten will. Erstere RE benennt die ganze Zeile, vom Beginn bis zum newline, zweiteere benennt nur das Zeichen '#' am Anfang der Zeile.

Eine Gleitkommazahl wie 3.675E-15 kann man mit [[[:digit:]]+\.[[:digit:]]\*([eE][+-]?[[:digit:]]+)? beschreiben. Zu beachten ist der Backslash '\' vor dem Punkt, da jener eine Sonderbedeutung hat, die mit dem vorangestellten '\' unterbunden wird. Leider hat diese Beschreibung einen Nachteil: sie lässt zwar die Zahl 1. durch, verbietet aber die Zahl .1, also noch einmal: (([[:digit:]]+\.[[:digit:]]\*)|(\.[[:digit:]]+))([eE][+-]?[[:digit:]]+)? Wie man sieht, werden REs schnell unübersichtlich. Das Chaos wird in

Verbindung mit **sed** und **bash** perfekt, da sich noch viele lustige `\` und `'` hinzugesellen werden. Dazu aber später.

## Grenzen von REs

Mit REs lassen sich nicht alle Zeichenketten beschreiben. Es ist zum Beispiel unmöglich ein System von balancierten Klammern zu beschreiben, auch ist die Menge  $\{w^c w \mid w \text{ ist ein String bestehend aus 'a's und 'b's}\}$  als RE nicht auszudrücken. Mehr zu REs kann man im 'Drachenbuch', *Compilers - Principles, Techniques and Tools* von Aho, Sethi und Ullman nachlesen.

## sed und REs

**sed** verwendet *Basic Regular Expressions*, eine Art Untermenge der oben vorgestellten Erweiterten Regulären Ausdrücke. Einige Unterschiede sind:

- Die Quantifikatoren `|`, `+` und `?` sind normale Zeichen, und es gibt keine äquivalenten Operatoren dafür. GNU-**sed** kennt diese Operatoren, wenn sie durch einen vorangestellten Backslash "escaped" werden
- Die geschwungenen Klammern sind normale Zeichen, und müssen mit Backslashes "escaped" werden, werden also als `\{` und `\}` geschrieben. Das selbe gilt für runde Klammern; die Zeichen, die durch `\(` und `\)` eingeschlossen werden, können später mit `\1` usw. dereferenziert werden
- `^` ist ein normales Zeichen, wenn es nicht am Beginn einer Zeile oder eines Klammerausdrucks steht
- `$` ist ein normales Zeichen, wenn es nicht am Ende einer Zeile oder eines Klammerausdrucks steht
- `*` ist ein normales Zeichen am Beginn einer Zeile oder eines Klammerausdrucks

## Erste Schritte mit sed

### Ein- und Ausgabe

**sed** liest von *stdin* und schreibt auf *stdout*. Man kann aber als letzten Parameter auf der Kommandozeile einen (oder mehrere) Dateinamen angeben, von dem die Eingabe gelesen werden soll. Weiters kann man sich der *Umleiteoperatoren* der Shell bedienen (`>`, `<`, `|`). Die drei folgenden Zeilen liefern das selbe Ergebnis:

```
sed -n -e '/root/p' /etc/passwd
sed -n -e '/root/p' < /etc/passwd
cat /etc/passwd | sed -n -e '/root/p'
```

Noch eine kleine Besserwisserei meinerseits, die mit **sed** eigentlich nichts zu tun hat, sondern mit Shell-Scripting. Von den beiden Zeilen

```
programm 2>&1 >file
```

und

```
programm >file 2>&1
```

ist die zweite Version vorzuziehen, da die erste Kommandozeile *stderr* auf den alte *stdout* setzt, und erst anschließend *stdout* auf *filename* umlenkt; *stderr* wird also i.d.R. nicht nach *filename* umgelenkt werden.

Die meisten UNIX-Kommandos lassen sich als *Filter* einsetzen. Filter werden dazu verwendet um einen Stream durch mehrere mit *pipes* (`|`) verkettete Programme zu jagen, jedes davon verändert den

Stream nach vorgegebenen Regeln. Auf diese Weise lassen sich in Verwendung von verschiedenen Filtern sehr komplexe Aufgaben erledigen.

## Kommandos

Das Programm

```
sed -e 'd' /etc/services
```

liefert erstmal gar nix.

Wie die Verarbeitung einer Zeile zu erfolgen hat, wird in einem *Script* oder *Programm*, festgelegt, das auf der Kommandozeile der Option `-e` folgen muss. Ein **sed**-Script enthält mindestens ein Kommando (in diesem Fall `d` für *delete*). Die Zeile wird also in den *pattern buffer* geladen, welcher dann nach den angegebenen Regeln bearbeitet (in diesem Fall gelöscht) und anschließend ausgegeben wird - die bearbeitete Datei wird dabei nicht verändert. Diese Schritte werden für jede Zeile wiederholt, bis zum Dateiende.

Bitte die beiden Anführungszeichen `' '` um das **sed**-Script `'d'` beachten, die eine Re-Interpretation der dazwischen liegenden Zeichen seitens der *Shell* verhindern. Diese Anführungszeichen sollten immer verwendet werden, da man sich dadurch viele unerwartete Reaktionen des Scriptes ersparen kann.

Vor manche Kommandos kann man eine *Adresse* stellen, die angibt welche Zeilen mit dem betreffenden Kommando zu bearbeiten sind. Somit kann man Kommandos selektiv auf bestimmte Zeilen (oder wie wir später sehen werden, auf bestimmte Zeichenketten) anwenden. Mehr zu Adressen im nächsten Kapitel.

## Adressen

Adressen können fixe Zeilen in einer Datei sein oder ganze Bereiche, oder aber Zeilen die auf einen bestimmten Reguläre Ausdruck passen.

```
sed -e '1d' /etc/services
```

Hier wird das Kommando `'d'` auf die Zeile mit der Adresse `'1'` angewendet. Der Effekt des Programms ist der, dass die erste Zeile von `/etc/services` in den *pattern buffer* geladen, dieser dann gelöscht und anschließend der leere *pattern buffer* (also nichts) ausgegeben wird. Alle anderen Zeilen werden in den *pattern buffer* geladen, der nicht bearbeitet wird da die Adresse nicht auf die Zeile zutrifft und anschließend wird der *pattern buffer* nach *stdout* geschrieben mit dem Resultat dass in der Ausgabe die erste Zeile fehlt.

Man kann auch Adress-Bereiche angeben wie in

```
sed -e '1,10d' /etc/services
```

was die ersten 10 Zeilen löscht oder man kann jede n-te Zeile bearbeiten wie in

```
sed -e '10~2d' /etc/services
```

wo jede zweite Zeile, ausgehend von der 10. Zeile gelöscht wird. Letzteres ist eine GNU-Erweiterung von **sed**; dort wo Portabilität auf andere Umgebungen wichtig ist, ist diese Adresse zu vermeiden.

Manchmal ist es interessant, nur solche Zeilen einer Konfigurationsdatei anzuzeigen, die nicht auskommentiert sind. Das kann mit folgender Zeile geschehen:

```
sed -e '/^#.*d' /etc/inetd
```

Die Adresse `/r/` wendet das nachfolgende **sed**-Kommando auf jede Zeile an, auf die der Reguläre Ausdruck `r` passt. Nur zur Erinnerung - die angegebene RE passt auf jede Zeile, welche "mit einem `#` beginnt und danach null oder mehr beliebige Zeichen enthält". Das ist aber nicht das was wir eigentlich wollten. Denn enthält die Datei eine leere Zeile, dann wird diese auch ausgegeben. Also müssen wir unsere Strategie ändern und z.B. nur jene Zeilen ausgeben, die mit einem Zeichen beginnen, das nicht `#` ist:

```
sed -e '/^[^#].*/p' /etc/inetd
```

Das ist neu: das Kommando `p`, das für *print* steht, also den *pattern buffer* ausgeben. Die Ausgabe ist aber alles Andere als erwartet: jede Zeile wird ausgegeben, die erwünschten Zeilen sogar zwei mal. Was ist passiert? Noch einmal müssen wir die Funktionsweise von **sed** durchkauen: Zeile einlesen, wenn die Adresse passt, dann *pattern buffer* bearbeiten (in unserem Falle ausgeben), dann *pattern buffer* ausgeben. Wir müssen also den letzten Schritt unterbinden; das geht mit der Option `-n` (portabel!) oder `--quiet` oder `--silent`, je nach Geschmack. Das richtige Programm schaut nun so aus:

```
sed -n -e '/^[^#].*/p' /etc/inetd
```

Wir haben gesehen, dass mit der Adresse `n, m` die `n`-te bis `m`-te Zeile bearbeitet wird - das geht auch mit REs: die Adresse `/BEGIN/, /END/` selektiert alle Zeilen ab der ersten Zeile, auf die die RE `BEGIN` passt bis zu der Zeile auf die die RE `END` passt oder bis zum Dateiende, je nach dem was früher kommt. Wird `BEGIN` nicht gefunden, dann wird keine Zeile bearbeitet. Es ist oft so, dass man beim Compilieren eines umfangreichen Projektes regelrecht von Fehlermeldungen und Warnungen erschlagen wird. Das ist ein Job für **sed**: das folgende Beispiel liefert nur jene Ausgaben des **gcc** die zwischen der ersten Warnung und der ersten Fehlermeldung liegen.

```
gcc sourcefile.c 2>&1 | sed -n -e '/warning:/,/error:/p'
```

Und wenn `n, m` gilt und `/BEGIN/, /END/`, warum nicht auch eine Kombination davon? Ein `/BEGIN/, m` heißt ab der Zeile, auf welche die RE `BEGIN` passt bis zur `m`-ten Zeile usw.

Das folgende Beispiel ist wohl das kürzeste sinnvolle Script in **sed**, das es gibt. Es gibt die Zeilenanzahl der bearbeiteten Datei aus (`wc -l`):

```
sed -n -e '$='
```

Das *Dollar*Zeichen '\$' ist in diesem Fall nicht das Zeilenende einer RE (es fehlen nämlich die `/`), sondern ist eine Adresse von **sed**, die die letzte Zeile des letzten Inputfiles benennt, und das Kommando '=' gibt die aktuelle Zeilennummer *vor* dem Output aus.

Ein Rufezeichen '!' nach einer Adresse kehrt diese in ihr Gegenteil um. Die Adresse `n, m!` trifft auf alle Zeilen bis auf den Block von der `n`-ten bis zur `m`-ten Zeile zu. `/awk/!` selektiert alle Zeilen die nicht die Zeichenkette 'awk' enthalten.

## Mehr Kommandos

Neben den Kommandos 'd' und 'p' die wir schon kennen gibt es noch eine Reihe anderer Kommandos, die ich aber nicht alle beschreiben kann. Hat man erst einmal die Syntax eines **sed** Programms verstanden, findet man sich leicht in der `man/info-page` zurecht und man kann sie dort nachschlagen.

Ein einfaches Kommando ist `q`, das das Programm abbricht. Ob der *pattern buffer* noch geschrieben wird hängt davon ab, ob die Option `-n` angegeben wurde oder nicht. Als Beispiel folgen zwei funktionsmäßig äquivalente Emulationen des UNIX-Befehls **head**, wobei die zweite Lösung effizienter ist, da sie nur die ersten 10 Zeilen bearbeiten muss.

```
sed -n -e '1,10p'  
sed -e '10q'
```

Weiters wird auch das Kommentarzeichen '#' als Kommando bezeichnet. Es verbirgt einfach alle nachfolgenden Zeichen im *Script* bis zum Ende der Zeile. Das ist nützlich in *Scripten*, die in Files geschrieben wurden und die an trickreichen Stellen ein paar erklärende Worte verlangen.

Ein wichtiges Kommando ist `'s/r/rep/flag'`. Hierbei wird diejenige Portion im *pattern buffer*, auf welche die RE 'r' passt durch die Zeichenkette 'rep' ersetzt und zwar in der Modalität, die mit dem *flag* bestimmt wird. Ein 'd' ersetzt das erste Muster und fängt dann einen neuen Zyklus an. Das Flag 'g' ersetzt *alle* Muster in einer Zeile, eine Nummer 'n' veranlasst `sed`, das `n`-te übereinstimmende Muster zu ersetzen. Mit dem Einzeiler

```
sed -e '/ich/s/$1500/$3000/g' Gehaltsliste.dat
```

kann man ein bisschen träumen. (Bei den Träumen wird es wohl bleiben, denn **sed** verändert die Datei nicht!) Wer jetzt denkt die Ausgabe mittels Ausgabeumleitung '>' wieder auf die Eingabedatei umzuleiten, der wird sich schön wundern: die Datei ist dann nämlich leer. Der richtige wenn auch umständliche Weg ist die Ausgabe in eine temporäre Datei umzuleiten und diese dann auf den Namen der Quelldatei umzubenennen. Noch nicht verstanden was das vorherige Beispiel gemacht hat? In der Zeile, die 'ich' enthält wird das Gehalt von \$1500 auf \$3000 verdoppelt, alle anderen Zeilen werden unverändert ausgegeben.

Die folgende Zeile lässt die Ausgabe des Shell-Kommandos **ls** hingegen sehr '1337' aussehen:

```
ls -l | sed -e 's/o/0/g' | sed -e 's/l/1/g' | sed -e 's/e/3/g'
```

Das ist als Kommandozeile ein wenig lang - könnte man nicht... Ja man kann das alles kompakter schreiben, indem man mehrere **sed**-Kommandos durch Strichpunkte trennt.

```
ls -l | sed -e 's/o/0/g;s/l/1/g;s/e/3/g'
```

Will man dem blinden Zorn des *Superusers* aus dem Wege gehen und eine Verhöhnung seiner Homedirectory vermeiden, muss man die Adresse `/ root$/!` den Kommandos voranstellen. Diese Adresse selektiert jede Zeile, die nicht mit 'root' endet. Um mehrere Kommandos auf eine Adresse zu binden, müssen diese gruppiert werden. Das geschieht mit den geschwungenen Klammern `{}`. Wichtig: auch nach dem letzten Kommando muss ein Strichpunkt gesetzt werden.

```
ls -l | sed -e '/ root$/!{s/o/0/g;s/l/1/g;s/e/3/g;}'
```

Das folgende Script zeigt dazu ein Beispiel und kann dazu verwendet werden, 8 Leerzeichen in ein tab zu verwandeln.

```
sed -e 's/ \{8\}/^t/g'
```

wobei das `^t` ein tab-Zeichen symbolisieren soll. Alles schön und gut, nur ist die tab-Taste unter der Shell für das schöne Wort Kommandozeilenervollständigung reserviert, ein Tabulatorzeichen selber kann man nicht direkt eingeben. Der einfachste Weg dazu ist die Tastenkombination `^V^I` zu drücken, was für **CTRL-V CTRL-I** steht. Ein `^V` fügt das nachfolgende Zeichen ohne weitere Interpretation auf der Kommandozeile ein. Alternativ kann man also auch `^V<tab>` tippen. Mehr dazu in den info-pages zu *bash*, *tsh* oder *readline*, sowie bei Ihrem Arzt oder Apotheker.

Es sei noch einmal angemerkt, daß Basic Regular Expressions die Zeichen `+` und `?` nicht kennen. GNU **sed** führt dagegen `\+` und `\?`. Da die gewünschten Effekte oft leicht mit Standard-Bordmitteln von **sed** zu erreichen sind, empfiehlt es sich, diese Konstrukte selten oder gar nicht zu verwenden.

Ein weiteres nützliches Kommando ist `'y/SOURCE-CHARS/DEST-CHARS/'`, das alle Zeichen in *SOURCE-CHARS* in das entsprechende Zeichen in *DEST-CHARS* umwandelt. Unnützlich zu sagen, dass beide Char-listen die gleiche Anzahl von Zeichen enthalten müssen. Das folgende Script 'verschlüsselt' den Text mit der sogenannten 'rot-13' Methode: alle Buchstaben werden um 13 Zeichen verschoben - aus 'a' wird 'n', aus 'b' wird 'o' usw. (Der Einfachheit halber werden hier nur Kleinbuchstaben verändert)

```
sed -e 'y/abcdefghijklmnopqrstuvwxyz/nopqrstuvwxyzabcdefghijklmnopghijklm/'
```

was auch ein schönes Beispiel für eine unüberlegte Benutzung von **sed** ist. Den selben Effekt kann man mit **tr** sehr viel einfacher und weniger fehlerträchtig erreichen:

```
tr '[a-z]' '[n-za-m]'
```

Ein Nachtrag zu den geschwungenen Klammern `{}`: Aus der Sicht von **sed** ist die öffnende Klammer `{` ein Kommando, dem eine Adresse oder ein Adressbereich vorangestellt werden kann. Das lässt sich für einen Trick missbrauchen, denn wenn man die Kommandos `=`, `a`, `i`, oder `r` (erlauben höchstens eine Adresse; zur Bedeutung dieser Kommandos bitte die Dokumentation bemühen) auf einen Adressbereich anwenden will, kann man sie in geschwungene Klammern setzen. So ist z.B. `'1,9='` ein ungültiges Kommando, aber `'1,9{=;}'` ist nicht zu beanstanden. Der Effekt dieses

Programms ist dass die Zeilen von 1 bis 9 mit vorangestellten Zeilennummern ausgegeben werden, der Rest des Files wird unverändert wiedergegeben.

Weil es oft gebraucht wird, stelle ich noch Scripte zur Umwandlung von Dateien im DOS-Format (CR/LF) ins UNIX-Format (LF) und umgekehrt vor. Sie wurden der schon erwähnten sedfaq [<http://sed.sourceforge.net/sedfaq.html>] von Eric Pement entnommen.

```
# 3. Under UNIX: convert DOS newlines (CR/LF) to Unix format
sed 's/.$//' file      # assumes that all lines end with CR/LF
sed 's/^M$//' file    # in bash/tcsh, press Ctrl-V then Ctrl-M
# 4. Under DOS: convert Unix newlines (LF) to DOS format
C:\> sed 's/$//' file      # method 1
C:\> sed -n p file         # method 2
```

Eine Randbemerkung: Ist keine `-e`-Option angegeben, dann wird der erste Parameter, der keine Option ist als das auszuführende Programm genommen. Um Verwirrung zu vermeiden empfiehlt sich immer ein `-e` anzugeben. Einem Guru wie Herrn Pement sei es aber gestattet sich über diese Faustregel hinwegzusetzen.

UNIX wäre nicht UNIX, wenn es nicht unzählige andere Methoden dafür gäbe: beispielsweise die Programme **dos2unix** bzw. **unix2dos**, oder der Befehl **tr -d [^M] < inputfile > outputfile** um vom DOS- ins UNIX-Format zu konvertieren, oder **:set fileformat=dos** bzw. **:set fileformat=unix** unter **vim** oder...

## Ein paar interessantere Beispiele

Wer bis hierher gekommen ist, sollte wirklich verstanden haben was *Adressen* und was *Kommandos* sind. Das ist wichtig, denn ab jetzt werden diese in einem **sed**-Script hintereinandergelinkt, und das kann sonst schon für einige Verwirrung sorgen.

Hin und wieder trifft man in Scripten nicht die gewohnte Form `/r/` einer RE vor - die *Slashes* `/` scheinen zu fehlen. Das hat den Grund, dass es manchmal nötig ist in einer RE den *Slash* selber anzugeben. Damit dieser aber nicht fälschlicherweise interpretiert wird, muss er mit dem *Backslash* gequotet werden, also `\`. **sed** gibt einem die Möglichkeit, ein anderes Zeichen als den *Slash* als RE-Begrenzer zu verwenden. Man kann also `/\bin\ls/` oder beispielsweise `\@/bin/ls@` verwenden. In gleicher Weise kann das mit dem `s-` oder `y-`Kommando geschehen: `'s//'` ist gleichwertig zu `'s@@'` Hat man deshalb nicht genau verstanden, was Adresse was Kommando und was RE ist, kommt man da leicht ins Schleudern.

## Probleme mit REs

Reguläre Ausdrücke finden immer den längsten passenden String. Das kann manchmal unerwünscht sein. Will man zum Beispiel eine HTML-Seite in Text umwandeln, dann könnte man in Versuchung kommen folgendes Script zu verwenden:

```
sed -e 's/<.*>//g' text.html
```

Das liefert aber nicht den gewünschten Effekt, denn eine Zeile

```
Das <b>ist</b> ein <i>Beispiel</i>.
```

wird zu

```
Das.
```

verküppelt. Man muss also nur jene Zeichen bis zum *ersten* `'>'` löschen:

```
sed -e 's/<[^>]*>//g' text.html
```

Muss man einen Text nicht bis zum ersten Vorkommen eines Zeichens sondern einer Zeichenkette bearbeiten, wird die RE ein bisschen komplizierter. Im Kapitel mit den Beispielen findet sich dazu ein Lösungsansatz (Löschen von Kommentaren).

## Selektives Ersetzen

Das `s///` Kommando kann nicht nur fixe Strings einsetzen, sondern auch den gefundenen String oder Substrings davon. Der *Ampersand* `&` steht dabei für den gesamten gefundenen String.

In meiner Kindheit hatten wir die *elleff*-Sprache, unsere Geheimsprache, bei der man jeden Vokal (oder Gruppe von Vokalen) in einem Wort mit `<VOKAL>l<VOKAL>f<VOKAL>` ersetzen muss. Kompliziert? Da ist die **sed**-Schreibweise prägnanter:

```
sed -e 's/[aeiou][aeiou]*/&l&f&/g'
```

Die Mächtigen der Welt, als 'Bilifill Clilifintolofon' oder 'Boloforilifis Jelefelzilifin' ausgesprochen, gewinnen damit in meinen Augen sofort an Sympathie. Meine Hochachtung jedem, der ein *verellefftes* 'ukulele' aussprechen kann ohne es vom Bildschirm zu lesen.

Unter GNU-**sed** kann man folgende Zeile schreiben:

```
sed -e 's/[aeiou]\+/&l&f&/g'
```

Bitte den *Backslash* `\` vor dem Plus beachten, da dieses Zeichen - weil GNU-Erweiterung - zuerst als normaler Charakter angesehen wird und seine Bedeutung die er bei REs inne hat, erst durch den Backslash gewinnt. Gleiches gilt auch für das *Fragezeichen* (*Questionmark*) `'?`, nicht aber für den *Asterisken* `'*`.

Hier weise ich noch einmal auf die Grenzen von Regulären Ausdrücken hin. Es ist nicht möglich, die Rücktransformation aus der *elleff*-Sprache mit REs auszudrücken. Ein `[aeiou]l[aeiou]f[aeiou]` kann man wohl angeben, nicht aber die Bedingung dass alle drei Vokale gleich sein müssen. Ob dies hinreichend ist um die *elleff*-Sprache als sichere Verschlüsselungsmethode zu bezeichnen, müssen wohl findigere Kryptologen entscheiden.

Mit **sed** ist es auch möglich, Teile von Strings heraus zu picken um diese später zu verwenden. Diese Teile werden mit `'(` und `)'` markiert, und man kann auf diese Strings mit `'\1'`, `'\2'` usw. zugreifen. Nehmen wir einmal an wir hätten ein File, in dem verschiedene Namen eingetragen sind:

```
John Fitzgerald Kennedy
Franz Josef Strauss
Ernst Theodor Amadeus Hoffmann
Theo Lingen
```

die in die Form `<VORNAME> [ <INITIAL ZWEITER NAME> . ] <NACHNAME>` gebracht werden soll. Dazu muss man erst die Regionen definieren:

```
sed -e 's/^[^ ][^ ]* [[[:alpha:]]]* [^ ][^ ]*$//'
```

Nun gibt man um die gewünschten Zonen die Klammern und stellt sich das Ergebnis mit `'\1'` und `'\2'` und `'\3'` zusammen:

```
sed -e 's/\([^[^ ]*\)\ ([[[:alpha:]]\)]* \([^[^ ]*\)\$/\1 \2. \3/'
```

und voilà das Ergebnis:

```
John F. Kennedy
Franz J. Strauss
Ernst T. Hoffmann
Theo Lingen
```

Will man das Ergebnis noch in eine Adressdatenbank importieren, dann muss man einen Feldbezeichner vor die Namen setzen. Ein erster Versuch wäre der, das gleich in einem Rutsch mit dem Script

```
sed -e 's/\([^[^ ]*\)\ ([[[:alpha:]]\)]* \([^[^ ]*\)\$/name: \1 \2. \3/'
```

zu bewerkstelligen, das liefert aber genau da ein falsches Ergebnis, wenn der zweite Vorname fehlt.

```
name: John F. Kennedy
name: Franz J. Strauss
```

```
name: Ernst T. Hoffmann
Theo Lingen
```

Einem solchen nur teilweise formatierten Datenschwulst ist nur schwer beizukommen. Deshalb den Output ungetesteter Scripte immer zuerst auf eine temporäre Datei umleiten, diese auf Korrektheit prüfen und dann die Zieldatei ersetzen. Wie man die Namen nun richtig formatiert, wird im nächsten Kapitel beschrieben. Warum hat das Script aber nicht richtig gearbeitet? Damit die RE auf eine Zeile zutrifft, muss diese mindestens 3 Felder, durch Leerzeichen getrennt, enthalten. Das ist bei Herrn Lingen nicht der Fall, deshalb wird auch das Kommando nicht ausgeführt und der *pattern buffer* wird unberührt gelassen.

## Mehrere Kommandos

Ein **sed**-Script kann mehrere Kommandos enthalten, die nach einander abgearbeitet werden. Man trennt diese mit einem *Semicolon* (;) oder man gibt diese mit mehreren `-e` Optionen an der Kommandozeile an oder man speichert die Kommandos in eine Datei die dann mit der Option `-f` *datei* abgearbeitet wird.

Eine mögliche Lösung des obigen Problems benutzt zwei Kommandos: das erste kürzt den Namen, ein zweites setzt vor alle Zeilen den String *name:*.

```
sed -e 's/\([^ ]*\) \([[:alpha:]]*\).* \([^ ]*\)$/\1 \2. \3/' \
    -e 's/./name: &/'
```

oder man trennt die zwei Anweisungen durch einen Strichpunkt (;). Zu beachten ist in der zweiten Anweisung die RE `. *`; würde man nur einen Punkt schreiben, passte dieser Ausdruck auch auf leere Zeilen. Das wird mit zwei Punkten vermieden.

Dieser Einzeiler in eine Datei geschrieben schaut so aus:

```
s/\([^ ]*\) \([[:alpha:]]*\).* \([^ ]*\)$/\1 \2. \3/
s/./name: &/
```

Und wieder eine Bemerkung die nichts mit **sed** zu tun hat: Die *Shell* gibt einem die Möglichkeit Scripte wie normale Programme zu behandeln. Dazu muss man nur an den Anfang des Scriptes die Zeile `#!/pfad/zum/interpreter <eventuelle Optionen>` setzen und die Scriptdatei als ausführbar markieren. Wenn diese Datei nun gestartet wird, ruft die Shell den angegebenen Interpreter mit dem Scriptnamen als Parameter auf. Auf das vorhergehende Beispiel angewandt sieht das so aus:

```
#!/bin/sed -f
s/\([^ ]*\) \([[:alpha:]]*\).* \([^ ]*\)$/\1 \2. \3/
s/./name: &/
```

Die Option `-f` weist **sed** an, den nachfolgenden Dateinamen (den die Shell hinzufügt) als Script zu nehmen. Dieser Trick funktioniert nur mit Scriptsprachen, bei denen das Zeichen `#` einen Kommentar einleitet, da sonst auch die erste Zeile als Programmcode interpretiert wird.

## spaceballs

### Ergänzungen zum *pattern space*

**sed** kennt noch weitere Kommandos zur Manipulation des *pattern space*. Das Kommando `D` löscht den Inhalt des *pattern space* bis zum ersten newline. Ist darin anschließend noch Text enthalten, wird ein neuer Zyklus gestartet, *ohne* eine neue Input-Zeile einzulesen. Ist der *pattern space* degegen leer, beginnt ein normaler Zyklus. Das Kommando `N` hängt ein newline an den Inhalt des *pattern space*, liest eine neue Zeile ein, welche nach dem newline eingefügt wird. Kann keine neue Zeile mehr eingelesen werden (Dateiende) dann wird das Programm an dieser Stelle abgebrochen. Das Kommando `P` gibt den Inhalt des *pattern space* bis zum ersten newline aus.

Das Beispiel dazu löscht alle *konsekutiven* Leerzeilen in einer Datei. Ist am Dateianfang eine Leerzeile, so bleibt sie erhalten, am Dateiende werden alle Leerzeilen gelöscht.

```
sed -e '/^$/N;/\n$/D'
```

## Einmal *hold space* und zurück

Neben dem *pattern space*, in den die Zeile geladen und dort manipuliert wird, kennt **sed** noch den *hold space*, der zu Programmbeginn leer ist, aber durch verschiedene Befehle manipuliert werden kann. Der *hold space* wird hauptsächlich dann verwendet wenn man das Operationsfeld eines einzigen Kommandos auf mehrere Zeilen ausdehnen will oder sich Zeilen für später aufheben muss.

Das Kommando 'h' überschreibt den *hold space* mit dem Inhalt des *pattern space*; die umgekehrte Operation wird durch das Kommando 'g' erreicht. Es gibt auch groß geschriebene Versionen dieser Kommandos, welche den Zielspace nicht überschreiben, sondern daran ein newline gefolgt vom Inhalt des Quellspace anhängen.

Zur Verinnerlichung des Konzepts des *hold space* ein sehr einfaches Beispiel, in dem die erste Zeile zurückbehalten wird und erst nach der letzten Zeile geschrieben wird. Das Programm kopiert also die erste Zeile in den *hold space*, gibt alle anderen aus, und nach Erreichen des Dateiendes wird der Inhalt des *hold space* in den *pattern space* kopiert, der dann noch ausgegeben werden muss. Das und nichts anderes tut der folgende Einzeiler.

```
sed -n -e '1h;1!p;${g;p;}'
```

Das folgende Beispiel gibt alle Zeilen sofort aus, die nicht in einem '/begin/,/end/'-Block liegen, den Rest erst bei Dateiende. Im Hinblick auf ein **sed**-Programm heißt das, Zeilen im Block '/begin/,/end/' werden an den *hold space* angehängt. Zu beachten ist nur, dass der Befehl 'H' dem Inhalt des *hold space* zuerst ein newline und dann der *pattern space* anhängt. Deshalb muss man bei der Ausgabe das erste Zeichen (sicher ein newline) unterdrücken.

```
sed -n -e '/begin/,/end/H;/begin/,/end/!p;${g;s/^././;p;}'
```

Anzumerken ist hierbei noch dass **sed** den Inhalt des *pattern space* als eine Zeile ansieht, egal ob da noch ein oder mehrere newline enthalten sind. Aus diesem Grund unterdrückt das Kommando 's/^././' nicht alle Buchstaben nach einem newline, sondern wirklich nur das erste Zeichen im *hold space*.

Das Kommando 'G' hat folgenden Effekt: es wird an den *pattern space* ein newline und anschließend der Inhalt des *hold space* angehängt. Das kann man für die verschiedensten Zwecke ausnützen. Das Script

```
sed -e 'G'
```

fügt nach jeder Zeile ein Leerzeichen ein (der *hold space* ist ja leer). Mit **sed** kann man auch die Funktionsweise von **tac** (ein umgekehrtes **cat**; dreht die Reihenfolge der Zeilen um) nachbilden:

```
sed -n -e 'G;h;${p;}'
```

mit dem kleinen Schönheitsfehler dass am Ende eine Leerzeile zu viel ausgegeben wird - sie ist die Leerzeile, die in der ersten Zeile dem *pattern space* unnötigerweise angehängt wurde. Diesen Fehler beheben gleich beide folgenden Programme.

```
sed -n -e 'G;h;${s/././;p;}'
sed -n -e '1!G;h;${p;}'
```

Mit dem Kommando 'x' werden die Inhalte der beiden spaces ausgetauscht. Abschließend zu diesem Kapitel möchte ich ein längeres Beispiel (Danke an Ulf Bro) vorstellen, das umgebrochene Absätze in eine einzelne Zeile umwandelt:

```
# Zeilen, die nicht leer sind werden dem Hold-Raum angehängt
# Bei Leerzeilen wird der Inhalt des Hold-Raums in den
# Pattern-Raum verlagert. Der Hold-Raum wird entleert
# Erste Newline wird entfernt, die anderen in Leerzeichen
# umgewandelt
/^$/! H
/^$/ {
```

```

x
s/\n//
s/\n/ /g
p
}
# Letzte Zeile nicht vergessen
$ {
  g
  s/\n//
  s/\n/ /g
  p
}

```

## Sprünge (*branches*)

Dieses Kapitel ist für *Stirb langsam* sed-Programmierer (frei aus dem Englischen übersetzt) geschrieben.

## Sprungkommandos

Sprungziele (*labels*) werden durch einen Doppelpunkt, gefolgt vom Namen des Labels gekennzeichnet '`: label`' wobei *label* ein beliebiger Name sein kann. Einen *unbedingten Sprung* (es wird also immer gesprungen) kennzeichnet man mit '`b label`' (b für branch). Das Sprungziel *label* muss natürlich irgendwo im Script definiert sein. Wird kein Label angegeben, dann fängt unmittelbar der nächste Zyklus an. Das Kommando '`t label`' definiert einen *bedingten Sprung*. Gesprungen wird, wenn im aktuellen Zyklus eine erfolgreiche Substitution ('`s //`'-Befehl) durchgeführt werden konnte und außerdem seither kein '`t`'-Sprung durchgeführt wurde. Auch hier gilt, wenn das Sprungziel nicht angegeben wurde, beginnt ein neuer Zyklus.

Das sind alle Kommandos in diesem Zusammenhang. In diesem Sinne kann man **sed** als echten RISC-Editor bezeichnen (RISC = Reduced Instruction Set Computer).

Achtung bei der Verwendung von '`t`', die manches Kopfzerbrechen bereiten kann. Das Große Reformationskript soll dies verdeutlichen.

```

#!/bin/sed -f
s/foo/bar/g
s/Bayern/Bayern/g;t noconversion
s/Katholik/Protestant/g
s/kathol/luther/g
: noconversion
# und weiter gehts im Code

```

Außer der Tatsache, dass man mit einem '`/Bayern/{ . . . }`' besser bedient wäre, sollte der Sinn des Scriptes klar sein: In jenen Zeilen, in denen das Wort 'Bayern' nicht vorkommt, soll alles Katholische durch Protestantisches ersetzt werden. Das eigentlich nutzlose '`s/Bayern/Bayern/g`' stellt die Bedingung für den nachfolgenden Sprung dar. Diese Zeile alleine? Nein, denn das Kommando '`s/foo/bar/g`' kann genau so gut ausgeführt werden und den 2 Zeilen entfernten Sprung einleiten. Denn obwohl dies durch die eigenwillige Formatierung des Scriptes so aussieht, ist das '`t`' Kommando nicht exklusiv an das unmittelbar davor stehende Kommando gebunden. Um einen Seiteneffekt durch das '`foo-bar`' Kommando zu vermeiden sollte man es tunlichst irgendwo unterhalb des Sprungkommandos unterbringen oder wenn das nicht möglich ist, dann muss ein *dummy*-Sprung eingeführt werden.

```

#!/bin/sed -f
s/foo/bar/g
t dummy
: dummy
s/Bayern/Bayern/g;t noconversion
s/Katholik/Protestant/g
s/kathol/luther/g

```

```
: noconversion
# und weiter gehts im Code
```

## Andere Sprünge

Auch andere Befehle wie 'q', 'd' und 'D' (bei Dateiende auch 'n' und 'N') verändern den Programmfluss. Bedacht eingesetzt, kann man mit den Sprungbefehlen von **sed** ziemlich komplexe Programme schreiben. Unbedacht eingesetzt, kann man unnötigerweise noch viel komplexere Programme schreiben.

## Vermischtes

### Dateien

Mit **sed** kann man auch Dateien lesen und schreiben. Das geht mit den den Kommandos 'r *filename*' und 'w *filename*'. Beim Lesen wird die Datei nach dem gegenwärtigen Zyklus ausgegeben, oder wenn eine neue Zeile gelesen wird. Eine nicht vorhandene Datei wird als existent aber leer angesehen. Der 'w' legt eine neue Datei an oder *überschreibt* eine schon vorhandene Datei und füllt sie mit dem Inhalt des *pattern space*. Das Kommando 'w' kann auch als Flag zu 's / /' angegeben werden, wobei in die Datei geschrieben wurde, wenn eine Substitution erfolgen konnte.

Das folgende **sed**-Script ist ein Ersatz für den UNIX-Befehl '**tee *dateiname***',

```
sed -e 'w dateiname'
```

und wenn man es in der Form '**sed w*dateiname***' schreibt, nur um ein 'w' länger als die Version mit **tee**

Um den Einsatz vom Kommando 'r' zu demonstrieren möchte ich meinen *Tante Amalien*-Emulator vorstellen (der von Joseph Weizenbaums genialem ELIZA inspiriert ist). Meine Tante Amalie hat 3 Standardsätze die sie der Reihe nach verwendet. Diese sind 'Meinst du?', 'Früher war es besser - entschieden besser!' und 'Davon verstehst du nichts.'. Diese in die 3 Dateien *stdsatz1* - *stdsatz3* geschrieben ergeben mit dem (unportablen!) Script

```
#!/bin/sed -nf
1~3r stdsatz1
2~3r stdsatz2
3~3r stdsatz3
```

das interessante Gespräch

```
Schönes Wetter heute
Meinst du?
Ja natürlich. Schau doch raus!
Früher war es besser - entschieden besser!
Ich weiß nicht was du hast - schöner als so kann es ja nicht sein.
Davon verstehst du nichts.
Wie du meinst, Amalie.
Meinst du?
```

usw. Ich könnte mich stundenlang mit Amalie unterhalten. Für komplexere Charaktere wäre es denkbar, den Input nach bestimmten Mustern zu durchsuchen und dementsprechend zu reagieren. Dabei muss ja nicht jede Antwort in einem File gespeichert sein, man könnte ja den *pattern buffer* mit 's/^.\*\$/Antwort/p' füllen.

## Noch mehr Kommandos

**sed** kennt noch einige meines Erachtens recht exotische Befehle wie 'a', 'i' und 'c', welche einen gegebenen Text ausgeben oder 'l', welcher nicht druckbare Zeichen im *pattern space* in ihrer ANSI C-Notation ausgibt und noch derer mehr. Sollten diese gebraucht werden, dann gibt die info-page bereitwillig Auskunft dazu.

## To sed or not to sed?

Carlos Jorge G. duarte emuliert in seinem Tutorial sehr viele UNIX-Befehle mit **sed**. Das mag als Beispiel sehr interessant sein (und ich empfehle jedem, diese Programme anzuschauen und zu verstehen), ist aber in der Praxis zu kompliziert. **sed** eignet sich durch die kompakte Schreibweise seiner Scripte besonders gut für Einzeiler - wird das Problem größer, kann das Debuggen eines **sed**-Scriptes sehr nervenaufreibend sein. Eine sehr gute Alternative ist da **awk**, das eine ähnliche Syntax (`/regex/{action}`) unterstützt, darüber hinaus noch strukturierte Programme (mit Konstrukten wie `'if (expr) statement else statement'` oder `'while (expr) statement'` u.Ä.) erlaubt, Funktionen zur Mustersuche, eigene Typen zur Behandlung von Gleitkommazahlen und vieles mehr besitzt. Leider kennt es die sehr bequemen `'\1'` Referenzen nicht, die **sed** bietet. **awk**-Programme sind in der Regel 3-10 mal so groß wie **sed**-Scripte.

**Python** oder **perl** bemüht man für größere und komplexere Probleme. Damit kann man ausgewachsene Programme schreiben, die zur Laufzeit interpretiert werden, wodurch sie etwas langsamer als **sed**- oder **awk**-Scripte abgearbeitet werden. Die Scripte können ungefähr die 8-40-fache Größe eines äquivalenten **sed**-Scriptes haben. Das muss kein Nachteil sein, denn die verlorene Zeit die ein mittelmäßiges **Python**-Script zur Laufzeit gegenüber einem **sed**-Script verliert, ist meistens durch eine lange Fehlersuche während des Schreibens eines solchen mehr als wett gemacht.

*Keines* dieser Tools soll verwendet werden, wenn das Betriebssystem ein dediziertes Programm für das jeweilige Problem bereitstellt, da dieses meist schneller und sicherer arbeitet und obendrein noch einfacher zu bedienen ist. UNIX bietet eine Vielzahl solcher Helferchen, die aber meistens nur einseitig oder überhaupt nicht verwendet werden. Ein Blick in die man-pages zu **bash/tcsh, xargs, tr, test/** [, **cat/tac, cut, [ef]grep, uniq, sort, wc, tail, column/colrm, paste, look, hexdump, basename** und **Kumpanen, seq, luser, lart, whack, bosskill**, [...] kann dem geeigneten Leser nur zum Nutzen gereichen.

## Andere Programme mit sed-Kommandos

Das meistverwendetste Programm, in dem man **sed**-Anweisungen angeben kann, ist sicher **vi**. Um beispielsweise in einer Zeile "bla" nach "blupp" zu ändern, schreibt man

```
:s/bla/blubb/g
```

Eine interessantes Flag zum `s///`-Kommando ist `c`, welches bewirkt, dass man bei jeder potentiellen Substitution um Bestätigung gefragt wird. Ein Blick in die Dokumentation von **vi** (**vim**) zu weiteren **sed**-Kommandos lohnt sich.

**ed** ist ein vollständiger Editor, und älter als **sed**. **ed** ist die richtige Wahl, wenn man Texte inline-editieren will. Zu beachten ist, dass, im Gegensatz zu **sed**, die Datei *nicht* tumb einmal von oben nach unten durchgeackert und eine Anweisung damit auf alle Zeilen angewandt wird. Das Kommando

```
s/bla/blupp/g
```

wird in **ed** nur auf die aktuelle Zeile angewandt. Der gewünschte Effekt wird mit

```
%s/bla/blubb/g
```

erreicht. Ein anderes neues Kommando is `g/re/command`. Es wendet das Kommando auf alle jene Zeilen an, auf welche der Reguläre Ausdruck 're' passt. Ein Beispiel dafür ist `g/re/p`, welche alle Zeilen ausgibt, welche auf die RE passen. Eine geschichtliche Randnotiz: das UNIX-Kommando **grep** wurde bequemlichkeitshalber aus **ed** extrahiert, und der Name deutet an, was das Programm tut: global-RE-print.

## Ein paar Beispiele

Diese Sektion sollte mehr Beispiele enthalten. Ich bin ständig auf der Suche nach Beispielen, welche zum Verständnis von **sed** beitragen. Sollte der Leser Scripte wissen, die mit noch nicht vorgestellten

Tricks arbeiten, welche eines Kommentars bedürfen, einem ganze Mannjahre an Handarbeit ersparen, einfach nur schön sind oder irgend einen anderen AHA!-Effekt auszulösen imstande sind, so bitte ich darum, mir diese zu schicken. Sie werden mit Angabe des Autors hier veröffentlicht.

In diesem Kapitel werden die GNU-Erweiterungen scham- und vor allem kommentarlos verwendet, da sie die Lesbarkeit eines Scriptes sehr verbessern. Die Beispiele ließen sich auch ohne diese Erweiterungen beschreiben (und es wird empfohlen das auch zu tun, sobald ein Script auf andere Systeme übertragen werden könnte) das ginge aber auf Kosten der Verständlichkeit.

## Entfernen von Kommentaren

Die fiktiven Hochsprachen der 5. Generation K und K++ kennen zwei Arten von Kommentaren. Da wäre die nur in K++ verwendete Art 'kk.\*' (zwei einleitende 'k'), oder die in beiden Sprachen verwendete Form 'ko.\*ok', wobei sich ein solcher Kommentar über mehrere Zeilen erstrecken kann. Es soll ein **sed**-Script erstellt werden, das solche Kommentare (warum auch immer!) entfernen soll.

Das Entfernen von K-Komentaren benötigt ein paar Erklärungen. Nehmen wir einmal an, wir hätten den gesamten Kommentar im *pattern buffer*. Das Kommando 's/ko.\*ok//' geht aus dem Grund nicht, weil die ansonsten nützliche Eigenschaft von REs, den längsten zutreffenden String zu nehmen, hier unerwünscht ist. Sind zwei vollständige Kommentare in einer Zeile vorhanden, dann würde auch der unschuldigerweise *dazwischen* stehende Code entfernt werden.

Der zweite Anlauf ist ein 's/ko\([^o][^k]\)\*ok//g'. Achtung bei Konstrukten, welche Quantifikatoren ('\*', '\+', '\{ \}' ...) auf zwei oder mehrere Zeichen anwenden! Das Script arbeitet nur bei der Hälfte der Kommentare, und zwar bei jener Hälfte welche eine gerade Anzahl von Zeichen beinhaltet. Vom Zorn gepackt, schreibt man dann Sachen wie 's/ko\([^o]\*\([o][^k]\)\*\([^o]\*\)\*ok//g' welche zwar korrekt sind, aber völlig Praxisuntauglich. Ein solches Monsterprogramm kann allerhöchstens auf einem Großrechner vernünftig arbeiten. In diesen Situationen hilft es, eine verbale Beschreibung des Musters zu finden. Die könnte so aussehen: "Ein K-Kommentar beginnt mit 'ko', ihm folgen null oder mehr der oben beschriebenen Lesevorgängen [^o]\|o\+[^ok] plus einem Abschluss o\+k". Auf **sed**isch übersetzt bekommt man 'ko\([^o]\|o\+[^ok]\)\*o\+k'.

Nun muss nur noch sicher gestellt werden, dass nach einem 'ko'-Muster auch ein 'ok' im *pattern buffer* ist. Ist dem nicht so, dann sorgt der innere Loop (um das label `append`) dafür, dass ständig neue Zeilen mit dem 'N'-Befehl an den *pattern buffer* angehängt werden. Ein äußerer Loop (um das label `test`) sorgt dafür dass jene Zeilen richtig behandelt werden, in denen ein Kommentar geschlossen und anschließend ein neues mehrzeiligen Kommentar wieder aufgemacht wird.

```
#!/bin/sed -f

#lösche K++-Kommentare
/^[[:blank:]]*kk.*$/d
s/kk.*//g

#Wenn kein Kommentar gefunden wurde, dann nächster Zyklus.
: test
/ko/!b

#Hänge so lange neue Zeilen an den pattern buffer an,
#bis ein vollständiger Kommentar zusammengebracht wurde.
: append
/ok/!{N;b append;}

#lösche K-Kommentare die sich vollständig im pattern buffer befinden
s/ko\([^o]\|o\+[^ok]\)*o\+k//g

t test
```

Soll das Script für in den archaischen Sprachen C/C++ geschriebene Programme funktionieren, dann muss man es mit dem **sed**-(Pseudo-)Einzeiler "sed -e '/^#{/!{s/k/\//g;s/o/\\*/g;}' k-kommentar > c-kommentar" ummodellern.

Ein K-Kommentar beginnt mit 'ko', soweit ist alles klar. Anschließend folgt die längst mögliche Zeichenkette, die kein 'ok' enthält. Hier liegt der Hund begraben. Das abschließende 'ok' ist wieder trivial.

Nun zum Hund: Gesucht ist die längst mögliche Zeichenkette, auf welche der reguläre Ausdruck /ok/ nicht zutrifft. Es ist also gewissermaßen das Gegenteil von /ok/ gesucht.

Man kann sich nun in Anlehnung an ein prozedurales Vorgehen vorstellen, man suche das erste Auftreten von 'ok' innerhalb einer Zeichenkette. Dazu lese man immerwieder neue Zeichen von der Zeichenkette ein und untersuche die eingelesenen Zeichen.

Die gesuchte Teil-Zeichenkette besteht dann aus null oder mehreren "Lesevorgängen": *Längste Zeichenkette ohne 'ok' = \(\Lesevorgang\)^\**

Nach jedem Lesevorgang trifft man dann eine Fallunterscheidung, etwa von der Art: Es wurde kein 'o' eingelesen, es wurde ein 'o' aber kein 'k' eingelesen, etc. Man erhält so:

```
Lesevorgang = Fall_1 \| Fall_2 \| ... \| Fall_x
```

Wieder in Anlehnung an das prozedurale Vorgehen wird man zu Beginn der Überlegungen davon ausgehen, es werde pro Lesevorgang nur ein einzelnes Zeichen eingelesen. Ist dieses Zeichen dann von 'o' verschieden, trifft also der Ausdruck [^o] darauf zu, kann man mit dem nächsten Lesevorgang fortfahren, und man hat:

```
Fall_1 = [^o]
```

Ist das eingelesene Zeichen hingegen gleich 'o', dann könnte man in die Versuchung kommen zu prüfen, ob sich das nächste Zeichen von 'k' unterscheidet, in der Annahme, damit einen weiteren Fall eines Lesevorganges vollständig abgehandelt zu haben: Den Fall o[^k] nämlich! Träfe dieser reguläre Ausdruck auf die immerhin bereits zwei eingelesenen Zeichen zu, dann ginge man zum nächsten Lesevorgang über.

Aber hoppla! Das vorhin auf [^k] überprüfte Zeichen könnte ja wieder gleich 'o' sein, was zur Folge hätte, dass man beim nächsten Lesevorgang das zuerst eingelesene Zeichen auf 'k' überprüfen müsste. Solche Abhängigkeiten zwischen den Lesevorgängen sprengen aber das Konzept dieser Vorgehensweise und deuten darauf hin, dass der vorhergehende Lesevorgang im Prinzip weitergeführt werden muss.

Ist das erste Zeichen eines Lesevorganges also ein 'o', dann könnte diesem 'o' gleich eine ganze Folge weiterer 'o's folgen. Man muss also einen Ausdruck der Form o\+ einlesen, und zwar solange, bis man endlich ein Zeichen findet, das sich von 'o' unterscheidet. Ist dann dieses Zeichen nicht nur von 'o', sondern auch von 'k' verschieden, dann hat man insgesamt einen Ausdruck der Form o\+[^ok] eingelesen:

```
Fall_2 = o\+[^ok]
```

Von da aus kann man nun problemlos mit dem nächsten Lesevorgang fortfahren. Da das erste Zeichen aber nur entweder 'o' oder dann eben nicht 'o' sein kann, treten neben Fall\_1 und Fall\_2 keine weiteren Fälle mehr hinzu:

```
Lesevorgang = Fall_1 \| Fall_2 = [^o]\|o\+[^ok]
```

Bei jedem Lesevorgang findet man also ein einzelnes Zeichen [^o] oder einen Ausdruck o\+[^ok]. Erst wenn eine "o-Folge" mit dem Zeichen 'k' endet, wenn man also auf den "Abschluss" o\+k trifft, ist man am Ende.

Die null oder mehr Lesevorgänge [^o]\|o\+[^ok] liefern damit die längste Zeichenkette, die den Abschluss o\+k nicht enthalten. Obwohl damit das anfängliche Ziel, die längste Zeichenkette ohne 'ok' zu finden, knapp verfehlt worden ist, kann man mit diesen Überlegungen bequem den letztlich gesuchten Ausdruck eines K-Kommentars hinschreiben: *Ein K-Kommentar beginnt mit 'ko', ihm folgen null oder mehr der oben beschriebenen Lesevorgängen [^o]\|o\+[^ok] plus einem Abschluss o\+k:*

```
K-Kommentar = ko\([^o]\|o\+[^ok]\)*o\+k
```

Übrigens kann man mit dem bekannten Editor vim einen K-Kommentar einfach durch

```
K-Kommentar = ko.\{-}ok
```

definieren. Dabei bedeutet der Ausdruck `\{-}`, dass ähnlich wie bei `.*` eine beliebige Zeichenkette gesucht ist, aber nicht längste, sondern die kürzeste.

Vielen herzlichen Dank an Mathias Michaelis für dessen Beitrag zu diesem Tipp.

## elleff-Rücktransformation

Gelogen habe ich nicht, als ich behauptete, mit REs könne man keinen *elleff*-verschlüsselten Vokal beschreiben - das stimmt schon. Aber **sed** kann. Und das auf eine sehr trickreiche Weise. Zuerst das Script, kommentiert wird danach.

```
sed -e 's/\([aeiou]\+\)\1\1f\1/\1/g'
```

Wenn man diesen Kniff nicht schon einmal gesehen hat, muss man 2 (ich 3) mal hinschauen um zu verstehen, warum das funktioniert. Was mir an diesem Beispiel so gut gefällt ist, dass sobald die Klammer geschlossen wird, der Inhalt der eingeschlossenen Region schon in `\1` bereit steht und somit verwendet werden kann - auch innerhalb der RE. Die RE wird somit zur Laufzeit verändert. Das zeigt einerseits wie leistungsfähig **sed** ist und andererseits dass es auch manchmal recht knifflig sein kann die Scripte anderer zu verstehen.

Diesen Trick verdanke ich Carlos Duarte - ein weiterer Anreiz, in sein **sed** tutorial [[http://www.mds.mdh.se/~dat95abs/sed\\_tutorial.txt](http://www.mds.mdh.se/~dat95abs/sed_tutorial.txt)] hineinzuschauen.

## Verschachtelte Klammern

In manchen Fällen muss man auf das n-te Feld einer Zeile zugreifen. Die Quantifikatoren `\{n\}` sind dabei sehr nützlich. Will man beispielsweise das 3. Wort einer Zeile an den Zeilenanfang setzen, ist das Konstrukt

```
sed -e 's/^\([ ^]* *\)\{2\}\([ ^]* \)/\2\1/g'
```

nicht richtig, da die Referenz `\1` nur das zweite Wort enthält, nicht aber das erste und zweite. Abhilfe schafft da ein weiteres Klammernpaar, wie im folgenden Script:

```
sed -e 's/^\(\([ ^]* *\)\{2\}\)\([ ^]* \)/\3\1/g'
```

Dabei ist `\3` die Referenz auf das zweite Wort (`\2` referenziert das letzte Wort in `\1`). Hierbei ist man schon an die Grenzen der Verständlichkeit eines **sed**-Scriptes gegangen, und es ist zu überlegen, ob man mit anderen Programmen wie zum Beispiel **awk** nicht besser bedient ist.

Danke an Tillmann Bitterberg für diesen Tip.

## Kurzreferenz

Diese Kurzreferenz kann und will nicht die man- oder info-page zu **sed** ersetzen, sondern ist nur als eine Gedächtnisstütze gedacht.

## Adressen

**sed** kann mehrere Dateien abarbeiten, wenn man diese als Argumente übergibt. Diese Dateien werden als ein einziger input-stream behandelt. Achtung deshalb bei Zeilenangaben. Die Adresse `'1'` gibt deshalb nicht die erste Zeile in jeder Datei an, sondern die erste Zeile im input-stream; der nachfolgende Befehl wird also nur einmal ausgeführt.

**Tabelle 2. Adressen**

Adresse	Beschreibung
<i>nummer</i>	Selektiert die Zeile <i>nummer</i> im input-stream.
\$	Letzte Zeile im input-stream.
<i>/regex/</i>	Alle Zeilen, auf die <i>regex</i> passt. Alternativ kann auch " <i>%regex\%</i> " angegeben, wobei '%' ein beliebiges Zeichen ist.
<i>adresse1,adresse2</i>	Adressbereich: Alle Zeilen zwischen <i>adresse1</i> und <i>adresse2</i> , einschließlich der beiden Adressen.
<i>adresse!</i>	Alle Zeilen ausschließlich der in <i>adresse</i> angegebenen Zeilen.

## Kommandos

**Tabelle 3. Allgemeine Kommandos**

Kommando	Anzahl Adressen	Beschreibung
#	0	Kommentar, alle nachfolgenden Zeichen bis zum Newline werden nicht als Programmcode interpretiert.
{	0-2	Beginnt einen Block, der mehrere Kommandos beinhalten kann. Muss mit } abgeschlossen werden. Auf jedes Kommando innerhalb des Blocks muss ein semicolon ";" folgen.
=	0-1	Schreibt die aktuelle Position im Input-stream.
q	0-1	Beendet das Programm. Der Pattern Buffer wird nur dann geschrieben, wenn die Option -n nicht gesetzt wurde.
l	0-2	Schreibt den pattern buffer in C-Notation.
d	0-2	Löscht den pattern buffer und startet sofort einen neuen Zyklus.
p	0-2	Schreibt den pattern space nach stdout. (wird normalerweise nur in Verbindung mit der Option -n verwendet.)
n	0-2	Schreibt den pattern space (wenn -n nicht gesetzt ist) und ersetze den pattern space mit der nächsten Zeile. Wenn keine Zeile mehr zu lesen ist, beende das Programm.
<i>s/regex/rpl/flg</i>	0-2	Ersetzt <i>regex</i> durch <i>rpl</i> . Null oder mehrere <i>flg</i> geben an, wie das geschehen soll: 'g' ersetzt alle Zeichenketten in einer Zeile, auf die <i>regex</i> passt, 'p' führt eine Substitution durch und schreibt den pattern buffer, 'nummer' ersetzt nur die <i>nummer</i> -te Fundstelle.
<i>y/src/rpc/flg</i>	0-2	Ersetzt jedes Zeichen im pattern buffer, das in <i>src</i> vorkommt, mit dem entsprechenden Zeichen in <i>rpc</i> .

**Tabelle 4. Sprung-Kommandos**

Kommando	Anzahl Adressen	Beschreibung
<i>: label</i>	0	definiert das Sprungziel <i>label</i> . Siehe Kommandos b oder t, wie man Labels anspringt.
b <i>label</i>	0-2	Branch; unbedingter Sprung.
t <i>label</i>	0-2	Bedingter Sprung. Es wird zum <i>label</i> gesprungen, wenn auf den aktuellen pattern buffer eine <i>s///-</i> oder <i>y///-</i> Substitution erfolgt wurde.

**Tabelle 5. Kommandos im Zusammenhang mit dem Hold-buffer**

Kommando	Anzahl Adressen	Beschreibung
D	0-2	Löscht den Text im <i>pattern-space</i> bis zum ersten newline. Ist noch Text im <i>pattern-space</i> enthalten, starte einen Zyklus mit diesem Text, ansonsten starte einen normalen Zyklus.
N	0-2	Hängt eine newline an den <i>pattern-space</i> an, gefolgt von der nächsten Zeile des inputs. Ist das Ende der Datei erreicht, wird das Programm abgebrochen, ohne weitere Befehle abzuarbeiten.
P	0-2	Gibt den <i>pattern-space</i> bis zum ersten newline aus.
h	0-2	Ersetzt den Inhalt des <i>hold-space</i> mit dem des <i>pattern-space</i> .
H	0-2	Hängt ein newline an den <i>hold-space</i> , gefolgt vom Inhalt des <i>pattern-space</i> an.
g	0-2	Ersetzt den Inhalt des <i>pattern-space</i> mit dem des <i>hold-space</i> .
G	0-2	Hängt ein newline an den <i>pattern-space</i> , gefolgt vom Inhalt des <i>hold-space</i> an.
x	0-2	Tauscht den Inhalt von <i>pattern-space</i> und <i>hold-space</i> aus.

Daneben gibt es noch weitere Befehle wie a, i, c, r, w. Der Leser sei diesbezüglich mit einem freundlichen RTFM auf die man-page verwiesen.

## Versionsgeschichte

### Versionsgeschichte

Version 1.4	11. August 2008	thp
Einige Rechtschreibkorrekturen. Vielen Dank an Constantin Hagemeier. Umstellung von db2pdf auf fop.		
Version 1.3	13. April 2008	thp
Viele Rechtschreibkorrekturen. Vielen Dank an Kate (KDE Advanced Text Editor).		
Version 1.2	03. März 2008	thp
Neue Lizenz: Creative Commons Attribution-Share Alike 3.0 Unported.		
Version 1.1	19. Juli 2007	thp
Kleinere Verbesserungen und Rechtschreibkorrekturen. Vielen Dank an Alexander Kriegisch.		
Version 1.0	20. Februar 2007	thp
Kleinere Verbesserungen und Rechtschreibkorrekturen.		
Version 0.9	03. Mai 2005	thp
Verbesserungen von Mathias Michaelis zu den K-Kommentaren.		
Version 0.8	17. März 2004	thp
Umstieg auf xml.		
Version 0.7	09. Januar 2003	thp
Bessere Unterscheidung Basic/Erweiterte Reguläre Ausdrücke; erster Versuch, GNU-ismen aus den Scripten zu verbannen.		
Version 0.6	10. November 2002	thp
Beispiel für den x-Befehl in der Space-ball Sektion eingefügt. Danke an Ulf Bro.		
Version 0.5	6. September 2002	thp
History als Kapitel angelegt, wie es die FDL verlangt; Detailänderungen.		
Version 0.4	2. März 2002	thp
Detailverbesserungen und Kurzreferenz.		
Version 0.3	11. September 2001	thp
Gründlich überarbeitet und neu strukturiert.		
Version 0.2	September 2001	thp
wurde nie freigegeben.		
Version 0.1	7. September 2001	thp

Beginn der Arbeit am Tutorium.