

sed

sed is a standard Unix utility which is a "stream editor", allowing transformations of its input. It reads the stream input line by line, performs command/script actions (if allowed on that line), and outputs the *modified* line.

Typical uses of sed:

- Execute repetitive edits to one or more files.
- Convert data
- Select lines containing certain data

The actions to be performed can be in either the command line directly or can be a script inside a *scriptFile*:

```
sed options 'command' file1 ...  
sed options -f scriptFile file1 ...
```

sed processing steps for each input line:

1. Place the line in a buffer (pattern space)
2. Execute the command or commands on the buffer
3. Output the buffer to stdout

Create a SedExamples folder. When logged into a fox server, please cd to the **/usr/local/courses/ssilvestro/cs3423/sed** directory and copy all the files to your sed directory.

```
cs1713p0.c  
cs1713p0v2.c  
file1  
file3  
file4  
file5  
file6  
Linux  
fileSSN  
inventory.txt
```

Three types of addressing modes in sed:

- 0 – no address → operates on all lines of input
- 1 – one address → operates on all matching lines
- 2 – address range → operates on all lines within the range

Example 1: Only output to stdout all lines in cs1713p0.c that contain "printf"

```
$ sed -n '/printf/p' cs1713p0.c  
// about the safety of scanf and printf  
printf("%-10s %-20s %10s %10s %10s %10s\n"  
        printf("invalid input when reading student data, only %d valid  
values. \n"  
        printf("\tdata is %s\n", szInputBuffer);  
        printf("%-10s %-20s %10.2f %10.2f %10.2f %10.2f\n"
```

- '/printf/p' means to print lines containing "printf"
- The -n means to **not** output lines normally. When used with the "p", only the results of the "p" are returned.

Example 2: Create a new file2 from file1 replacing "Linux" with "Unix".

```
$ sed -i.bak 's/Unix/Linux/g' *.tex
```

- **s** means **substitute**. It is followed by a match pattern and a replacement value.
- **g** means **globally** change each occurrence on the line (this doesn't mean globally in the file)
- every line is sent to file2 including both modified and unmodified lines.

Example 3: Modify file3 removing carriage returns. For safety, make a backup of file3 in file3.sav.

```
$ sed -i.sav 's/\r/\n/g' file3
```

- -i means to edit the file "in place" which means to modify the file. It creates a backup copy of "file3" using the specified suffix ".sav". This option is very useful if you need to modify multiple files.
- replaces each carriage return with an empty string (i.e., deletes the carriage return)

Example 4: Modify file4 and file5 replacing only one occurrence of "cat" with "dog" per line for lines 1 thru 3. Save the old files with a ".cat" suffix

```
$ sed -i.cat '1,3 s/cat/dog/' file4 file5
```

- 1,3 is the range of lines to apply the substitution
- The backup copies are named file4.cat and file5.cat

Example 5: Delete lines in file6 that begin with "#", producing file6.new

```
$ sed '/^#/d' file6 > file6.new
```

- The '/^#/d' is a pattern for deleting lines.
- "^" specifies that the pattern must be at the beginning of the line.

	<ul style="list-style-type: none"> The "d" means to delete. <p>Note: the characters after the ending slash/delimiter are called flags.</p>
<p>Using a Script File</p> <p>As stated above, the -f switch is used to specify a script file which allows for more complex capabilities.</p> <p>Notice that example #1 printed lines which had "printf" in "/" comments. We can remove those lines.</p> <p>Our script will:</p> <ul style="list-style-type: none"> delete lines that begin with any number of spaces (including zero) followed by "/". Any number of spaces is specified by " *". print lines that contain "printf" 	<p>Example 6: Only output to stdout the lines in cs1713p0.c that contain "printf", but don't include lines that begin (not necessarily in column 1) with "/".</p> <p># Create a script file named "example6"</p> <pre>\$ cat >example6 /^ *\\//d /printf/p</pre> <p>CTRL-D</p> <ul style="list-style-type: none"> The first action gave a pattern that begins at the beginning of the line, followed by any number of spaces, and then followed by two slashes which had to be escaped using backslashes since slashes are our delimiters for patterns. If the pattern in the first action is matched, the line is deleted. The second action is only executed after the first action. If the line is deleted, the second action is not executed. The second action prints lines containing "printf" <pre>\$ sed -n -f example6 cs1713p0.c printf("%-10s %-20s %10s %10s %10s %10s\n" printf("invalid input when reading student data, only %d valid values. \n" printf("\tdata is %s\n", szInputBuffer); printf("%-10s %-20s %10.2f %10.2f %10.2f %10.2f\n" # Also try it on cs1713p0v2.c \$ sed -n -f example6 cs1713p0v2.c // about the safety of scanf and printf printf("%-10s %-20s %10s %10s %10s %10s\n" printf("invalid input when reading student data, only %d valid values. \n" printf("\tdata is %s\n", szInputBuffer); printf("%-10s %-20s %10.2f %10.2f %10.2f %10.2f\n"</pre> <p>What problems do you see with that output?</p> <p>??.</p>
<p>Example 6 only handled leading spaces. What if the "/" was preceded by tabs and/or spaces?</p> <p>We need a pattern of any number of spaces or tabs in any order. This is done by "[\t]*".</p> <p>Note that in non-GNU versions of sed, you may have to enter a tab character instead of "\t".</p>	<p>Example 7: Only output to stdout the lines in cs1713p0.c that contain "printf", but don't include lines that begin (not necessarily in column 1) with "/". The "/" could be preceded by a combination of spaces and tabs.</p> <p># Create a script file named "example7"</p> <pre>\$ cat >example7 /^[\t]*\\//d /printf/p</pre> <p>CTRL-D</p> <ul style="list-style-type: none"> The first action gave a pattern that begins at the beginning of the line, followed by any number of spaces or tabs, and then followed by two

	<p>slashes which had to be escaped using backslashes since slashes are our delimiters for patterns.</p> <ul style="list-style-type: none"> • If the pattern in the first action is matched, the line is deleted. • The second action is only executed after the first action. If the line is deleted, the second action is not executed. • The second action prints lines containing "printf" <pre>\$ sed -n -f example7 cs1713p0v2.c printf("%-10s %-20s %10s %10s %10s %10s\n" printf("invalid input when reading student data, only %d valid values. \n" printf("\tdata is %s\n", szInputBuffer); printf("%-10s %-20s %10.2f %10.2f %10.2f %10.2f\n"</pre>
<p>Example 6 and example 7 didn't print the entire printf statements; instead, those examples just printed the line that contained "printf". How can we print the entire printf statement?</p> <p>We need a range of lines that begin with a "printf(" and end with its ");"</p>	<p>Example 8: Only output to stdout the lines in cs1713p0.c that contain actual "printf" statements including all lines until the ");".</p> <p># Create a script file named "example8"</p> <pre>\$ cat >example8 /^[\t]*\\\/\\\/d /printf(/, /);/p</pre> <ul style="list-style-type: none"> • The second action lists a range of lines that begin with a line containing "printf(" and ends with a line containing ");". Those are two different patterns which are in a range since there is a comma between them. <pre>\$ sed -n -f example8 cs1713p0.c printf("%-10s %-20s %10s %10s %10s %10s\n" , "ID", "Name", "Exam 1", "Exam 2", "Final", "Average"); printf("invalid input when reading student data, only %d valid values. \n" , iScanfCnt); printf("\tdata is %s\n", szInputBuffer); return ERR_BAD_INPUT; } dAverage = (student.dExam1 + student.dExam2 + student.dFinalExam) / 3; printf("%-10s %-20s %10.2f %10.2f %10.2f %10.2f\n" , student.szStudentIdNr , student.szStudentFullNm , student.dExam1 , student.dExam2 , student.dFinalExam , dAverage);</pre> <p>Notice that the result is incorrect. Why?</p> <p>??</p> <p>Answer: Two-address mode cannot start and stop on the same line; the second address is not checked until at least the next line has been read, and thus cannot apply to a single line.</p>
<p>We can fix the problem with example 8 as shown to the right.</p>	<p>Example 9: Only output to stdout the lines in cs1713p0.c that contain actual "printf" statements including all lines until the ");". This fixes the problem in #8</p> <p># Create a script file named "example9"</p> <pre>\$ cat >example9 /^[\t]*\\\/\\\/d</pre>

	<pre>/printf(.*);/p /printf(.*);/d /printf(/,/);/p</pre> <ul style="list-style-type: none"> The second action matches a statement that begins with "printf(" followed by any characters and then has a ");" all on the same line. It then prints it. The third action does the same pattern match as the second action, but it deletes the line; therefore, the fourth action is only executed when it doesn't match. <pre>\$ sed -n -f example9 cs1713p0.c printf("%-10s %-20s %10s %10s %10s %10s\n" , "ID", "Name", "Exam 1", "Exam 2", "Final", "Average"); printf("invalid input when reading student data, only %d valid values. \n" , iScanfCnt); printf("\tdata is %s\n", szInputBuffer); printf("%-10s %-20s %10.2f %10.2f %10.2f %10.2f\n" , student.szStudentIdNr , student.szStudentFullNm , student.dExam1 , student.dExam2 , student.dFinalExam , dAverage);</pre>
<p>Some important pattern symbols for sed regular expressions:</p> <p>Basic regular expressions (requires "-r" in sed or "-E" in grep to use natively):</p> <p>.</p> <p>matches any character including newline</p> <p>^</p> <p>the pattern that follows must begin at the first character</p> <p>\$</p> <p>the pattern that precedes the \$ must match at the end of the line</p> <p>*</p> <p>matches zero or more of the preceding character, group or bracketed list</p> <p>[list]</p> <p>matches one character to any of the characters listed within the brackets. For convenience, range abbreviations can be used (e.g., [a-f], [0-9])</p> <p>[^list]</p> <p>matches one character if it is not listed within the brackets.</p> <p>Extended regular expressions (requires "-r" in sed or "-E" in grep to use natively, without needing escape characters):</p> <p>\(regexp\)</p> <p>groups the inner <i>regexp</i> and allows it to be referenced later with \1, \2, etc</p> <p>+</p> <p>similar to *, but matches one or more</p> <p>{i}</p> <p>matches <i>i</i> occurrences of the preceding character, group or bracketed list</p> <p>{i, }</p> <p>matches <i>i</i> or more occurrences of the preceding character, group or bracketed list</p> <p>{i, j}</p> <p>matches between <i>i</i> and <i>j</i> occurrences of the preceding character, group or bracketed list</p> <p>?</p> <p>matches zero or one instance of the preceding;</p>	<p>Example patterns for sed regular expressions:</p> <p>/^The/</p> <p>Matches any line that begins with the word "The"</p> <p>/[a-z]\{3\}[0-9]\{3\}/</p> <p>Matches a lowercase abc123 ID.</p> <p>/[0-9]*\$/</p> <p>Matches any line that ends with zero or more digits.</p> <p>Why is that pattern not very useful? What would make it have to end in digits? /[0-9]\+\$/</p> <p>\(bark\)\+ /</p> <p>Matches one or more occurrences of "bark"</p> <p>/(.*)/</p> <p>Matches a left parenthesis followed by any number of any characters followed by a right parenthesis.</p> <p>What does this pattern mean?</p> <p>/[^abcd]\.\$/</p> <p>??</p> <p>What would be the pattern for matching lines that are empty or contain any number of spaces and/or tabs?</p> <p>??</p> <p>Sam's note: demonstrate different between these</p> <p>Use --color flag to illustrate matches / nature of match.</p> <p>\grep "[^abcd].\$" ~/lines</p> <p>\grep "[^abcd]\.\$" ~/lines ← is now a literal, not a quantifier!</p> <p>\grep "[a-z]?\$" ../lines</p>

<p>equivalent to <code>{0,1}</code> <code>pat1 pat2</code> alternation: matches either <code>pat1</code> or <code>pat2</code></p> <p>Symbols representing characters to match (e.g., <code>.</code>, <code>[list]</code>, etc) are called classes. Symbols which dictate the amount of a preceding class (e.g., <code>*</code>, <code>+</code>, etc) are called quantifiers.</p> <p>Groups can be referenced later with <code>"\1"</code>, <code>"\2"</code>, etc.</p> <p>To reference the characters <code>"."</code>, <code>"^"</code>, <code>"\$"</code>, <code>"*"</code>, <code>"["</code>, <code>"]"</code>, in a regular expression, it must be escaped with a backslash.</p> <p>To reference the characters <code>"+"</code>, <code>"{"</code>, <code>"}"</code>, <code>"("</code>, <code>")"</code>, just use each normally. Their special (extended) meaning requires the use of a backslash when referenced in a basic regular expression.</p>	<p><code>\grep "[a-z]\?\$" ../lines</code> ← is now a quantifier, not a literal!</p>
<p>Note that sed also supports extended regular expressions if the <code>-E</code> or <code>-r</code> switch is specified (BSD vs GNU). This changes some of the above symbols and adds some additional symbols.</p> <p><code>.</code> matches any character including newline</p> <p><code>^</code> the pattern that follows must begin at the first character</p> <p><code>\$</code> the pattern that precedes the <code>\$</code> must match at the end of the line</p> <p><code>*</code> matches zero or more of the preceding character, group or bracketed list</p> <p><code>[list]</code> matches one character to any of the characters listed within the brackets. For convenience, range abbreviations can be used (e.g., <code>[a-f]</code>, <code>[0-9]</code>)</p> <p><code>[^list]</code> matches one character if it is not listed within the brackets.</p> <p><code>(regexp)</code> groups the inner <i>regexp</i> and allows it to be referenced later with <code>\1</code>, <code>\2</code>, etc</p> <p><code>+</code> similar to <code>*</code>, but matches one or more</p> <p><code>{i}</code> matches <i>i</i> occurrences of the preceding character, group or bracketed list (see other variants above (e.g. <code>{i,}</code>, <code>{1,i}</code>))</p> <p><code>pat1 pat2</code> matches either <code>pat1</code> or <code>pat2</code></p> <p><code>?</code> matches zero or one instance of the preceding; equivalent to <code>{0,1}</code></p>	<p>Example patterns for sed extended regular expressions:</p> <p><code>/[a-z]{3}[0-9]{3}/</code> Matches a lowercase abc123 ID.</p> <p><code>/(bark)+/</code> Matches one or more occurrences of "bark"</p> <p><code>/\(.*\)/</code> Matches a left parenthesis followed by any number of any characters followed by a right parenthesis.</p> <p><code>/(abc) (xyz)/</code> Matches exactly "abc" or "xyz".</p> <p>Sam's note: Look at <code>~/sed/samexp</code></p> <p>Using sed extended regular expressions, what would be the pattern to match a phone any phone number like (210)456-1234, but with any digits?</p> <p>Purely <i>basic</i> regular expressions: <code>([0-9][0-9][0-9])[0-9][0-9][0-9]-[0-9][0-9][0-9]</code></p> <p>Extended regular expressions, but in <i>basic mode</i>: <code>([0-9]\{3\})[0-9]\{3\}-[0-9]\{4\}</code></p> <p>Purely extended regular expressions: <code>\([0-9]{3}\)[0-9]{3}-[0-9]{4}</code></p>

<p>To reference the characters ".", "^", "\$", "*", "[", "]", "+", "{", "}", "(", ")", " " in an extended regular expression, it must be escaped with a backslash.</p>	<p>Run: <code>~/sed/phonepats.bash phonenum</code></p>
<p>Exercise #1: For privacy reasons, we want to replace Social Security Number references in emails from customers with XXX-XX-XXXX when we archive the email.</p> <p>Assume the input file is fileSSN, produce an output file with SSNs replaced as described above.</p>	<p>Sed in extended mode: <code>sed -r -e 's/[0-9]{3}-[0-9]{2}-[0-9]{4}/XXX-XX-XXXX/g' ~/sed/ssnnum</code> Sed in regular mode, but using extended regex: <code>sed -e 's/[0-9]\{3\}-[0-9]\{2\}-[0-9]\{4\}/XXX-XX-XXXX/g' ~/sed/ssnnum</code> Sed in regular mode, using purely basic regex: <code>sed -e 's/[0-9][0-9][0-9]-[0-9][0-9]-[0-9][0-9][0-9][0-9]/XXX-XX-XXXX/g' ~/sed/ssnnum</code></p>
<p>Negation We can negate an address (including a regex pattern) by following the pattern with a bang.</p>	<p>Example 10: Delete all lines in the file that don't contain "linux" or "Linux" <code>\$ sed '/[Ll]inux/! d' Linux</code> Although Linux is well known today, Linux was not known was a high school English teacher, if she knew Linux.</p>
<p>Append and Insert a Append inserts one or more lines of text after the matched lines. i Inserts one or more lines of text before the matched lines.</p>	<p>Example 11: We will append the Line "hello there" to each of the first three lines of the data from file5. # create a script file named "app11" <code>\$ cat > app1</code> 1,3a hello there CTRL-D</p> <p><code>\$ sed -f app11 file5</code> The weather was getting very bad. It was raining cats and dogs. hello there It was so bad that I stepped in a poodle. hello there what a catastrophe. hello there I was so tired that I needed a cat nap.</p> <p>Example 12: Insert the two lines shown below before any lines containing "cat" or "dog" in file5. # create a script file name "ins12" <code>\$ cat > ins12</code> /(cat) (dog)/i\ Who is watching\ a cat or dog video? CTRL-D</p> <p><code>\$ sed -E -f ins12 file5</code> Who is watching a cat or dog video? The weather was getting very bad. It was raining cats and dogs. It was so bad that I stepped in a poodle. Who is watching</p>

	a cat or dog video? what a catastrophe. Who is watching a cat or dog video? I was so tired that I needed a cat nap.
Change c changes one or more lines replacing them with the specified text	Example 13: Replace lines 2-6 with "eunuchs" for the linux file \$ sed '2,6c eunuchs' Linux Although Linux is well known today, Linux was not known eunuchs a word that was a homonym, she was unnerved. Example 14: Completely replace lines containing "Linux" or "Unix" with simply "eunuchs" \$ sed -E '/(Linux) (Unix)/c eunuchs' Linux eunuchs outside of technical circles in the 1980s. I worked for a company that was hiring technical writers in the mid 1980s. An interviewer asked a candidate, who eunuchs Having never heard of the operating system, but knowing a word that was a homonym, she was unnerved.
Substitute We have already seen several examples for substitute (examples 2, 3, and 4 and exercise #1). <i>s/matchPattern/replaceValue/flags</i> <ul style="list-style-type: none"> <i>matchPattern</i> is the pattern to find in a line <i>replaceValue</i> is the value used as a replacement for the matched value <i>flags</i>: <ul style="list-style-type: none"> <i>g</i> - globally replacement all occurrences in the line <i>p</i> - print the contents of the pattern space if successful <i>w</i> - write pattern space to a file if successful <i>i</i> - this is a number; replace the <i>ith</i> occurrence only <p>Although our examples have shown slash ("/") as the delimiter between the <i>matchPattern</i>, <i>replaceValue</i>, and <i>flags</i>, you can specify a different delimiter after the "s". This can make it easier when you need slash in your pattern or replacement value. Substitute simply uses the first character after the "s" as the delimiter</p> <p>Instead of <code>s/\V.*\$//</code> you could use <code>s / .*\$ </code></p>	Example 15: Examine the file inventory.txt. It contains many inventory records. The values are separated by spaces except the last data value is a product name which can contain multiple spaces and follows the phrase "Name:". We want to keep the product ID and unit price. This means that we want to remove the second, third, and fifth logical columns. Note that product IDs end in three numbers. \$ cat >example15 ?? CTRL-D \$ sed -f example15 inventory.txt PPF001 9.95 SBB001 14.95 SBG002 14.95 BOM001 29.95 MCW001 12.45 TTP001 9.95 NHC001 9.95 SSX001 29.95 <u>Sam's note:</u> my answer in ~/sed/example15.bash sed -r -e 's/ [0-9]+ [0-9]+ ([0-9.]+) .*\$/ \1/' inventory.txt
Alternative to Scripts when needing multiple commands	Example 15-2: use -e command arguments to solve example 15. \$ sed -e ?? -e ?? \

<p>The -e (lowercase) command argument can be used to provide multiple edits to a file without requiring a script file.</p> <p>This is very useful when you need variable values for different sed script steps since sed doesn't provide an easy way to pass variables as parameters.</p>	<pre>? '-e ??' inventory.txt Answer: sed -r -e 's/ [0-9]+// ' -e 's/ [0-9]+// ' -e 's/ ([0-9.]+) .*\$/ \1/' inventory.txt PPF001 9.95 SBB001 14.95 SBG002 14.95 BOM001 29.95 MCW001 12.45 TTP001 9.95 NHC001 9.95 SSX001 29.95</pre>
<p>Next</p> <p>The next command can be used to skip lines without printing them.</p> <ul style="list-style-type: none"> • n – Skips to the next line of input and overwrites the current line. Therefore, the current line and all changes made to it will not be printed. • N – Reads the next line of input and appends it (with a newline character preceding) to the end of the current line. <p>sed uses newline characters (<code>\n</code>) to separate one line from another and when a line is copied into the pattern space the newline character is removed. This means that using the commands we've seen so far, we would be unable to remove newline characters from a file. One solution to this is the N command.</p> <p>Note that when printing a line of output sed prints a newline character provided there is at least one character in the buffer. Because of these restrictions sed is not the recommended tool for removing newline characters. Instead use 'tr'</p>	<p>Example 16: Remove newline characters from any line containing 'Unix' and replace them with a tab character</p> <pre>\$ cat >example16 /(Linux) (Unix)/ s/\n\t/ CTRL-D \$ sed -E -f example16 < file1</pre> <p>Although Unix is well known today, Unix was not known outside of technical circles in the 1980s. I worked for a company that was hiring technical writers in the mid 1980s. An interviewer asked a candidate, who was a high school English teacher, if she knew Unix. Having never heard of the operating system, but knowing a word that was a homonym, she was unnerved.</p> <p>The above script fails to find and replace the newline characters because they have already been removed by sed before pattern matching happens.</p> <pre>\$ cat >example16 /(Linux) (Unix)/N s/\n\t/ CTRL-D \$ sed -E -f example16 < file1 Although Unix is well known today, Unix was not known outside of technical circles in the 1980s. I worked for a company that was hiring technical writers in the mid 1980s. An interviewer asked a candidate, who was a high school English teacher, if she knew Unix. Having never heard of the operating system, but knowing a word that was a homonym, she was unnerved.</pre>
<p>Input Buffers</p> <p>The sed editor provides 2 buffers to use while editing a file. These are referred to as the pattern space and the hold space.</p>	<p>The pattern space and hold space buffers have several commands that can be used to manipulate them.</p>

<p>All commands we have used so far operate on the pattern space.</p> <p>Immediately after finishing all commands in the file and before looping back, sed outputs the contents of the pattern space followed by a newline unless the -n flag is set. This is why all lines print by default.</p> <p>The hold space is an additional buffer that can be used to hold a value across multiple lines. It remains unchanged until operated on by one of the commands on the right.</p>	<ul style="list-style-type: none"> • g – Copies contents of hold space into pattern space • G – Appends contents of hold space onto pattern space (Adding a newline between) • h – Copies contents of pattern space into hold space • H – Appends contents of pattern space into hold space (Adding a newline between) • x – Swaps values stored in pattern space and hold space
<p>Move a line</p> <p>Example 16 first uses the h command to store the contents of line containing the pattern 'last'. Next the d command is used to remove it from the pattern space so that it is not printed. Finally, on the last line of input the G command is used to append the line stored in the hold space onto the end of the pattern space before printing.</p> <p>What would happen if there were multiple lines in the file containing the pattern 'last'?</p> <p>??</p>	<p>Example 17: Examine the file hold_example.txt. It has a problem; the last line is in the wrong position and needs to be moved to the end. This can only be done if we are somehow able to save off the first line and only print it when we reach the last line.</p> <pre>\$ cat >example17 /last/h /last/d \$G CTRL-D \$ sed -f example17 hold_example.txt First line Second line here Here is the third line Almost at the final line. This should be the last line.</pre>
<p>Warning: some implementations of sed restrict the max size stored in both the pattern space and hold space to 4000 bytes. GNU sed has no max buffer size provided it can malloc more memory.</p>	