

Perl Short Course: Third Session

Duncan C. White (d.white@imperial.ac.uk)

Dept of Computing,
Imperial College London

December 2011

A better way to run Perl Programs

Well, on Unix we can:

- First, issue the Unix command:

```
chmod +x eg1
```

This makes the file executable.

- Second, edit **eg1** and add the following line at the top:

```
#!/usr/bin/perl
```

- This is a special line interpreted by Unix when it executes a non-machine code program.
- Unix executes the named program (the Perl interpreter) with the script **eg1** as a command line argument.
- Perl starts up, reads **eg1** and proceeds to run it - and then ignores the first line because it's a comment!
- Now, run **eg1** by **eg1** (if **.** is on your path), or **./eg1** if not.

In this third session, we'll go over more of Perl in detail, we'll look at:

- **arrays and lists**
- **hashes**
- **special variables** (@ARGV, \$_, %ENV) and
- **regular expressions.**

Aside: A better way to run Perl Programs

- We have seen that when we want to run a Perl program called **eg1**, we say: `perl eg1`.
- Wouldn't it be better if we could just type `eg1` to run our program?
- Then we could install our own Perl programs in a public place and let our friends run them - without them caring what language the programs are written in!

Arrays and Lists

- An array is an ordered collection of scalars (strings or numbers), declared via `my @array`, the `@` being compulsory.
- An array such as `@fred` is *not the same as* `$fred`. Perl keeps the namespaces of arrays and scalars separate.
- Array indices start at 0, like C, C++ and Java, but unlike many other languages.
- An array may be built up piece by piece:

```
my @fred;
$fred[0] = "hello";
$fred[1] = 7.1+$a;
$fred[2] = 17.3;
$fred[3] = $c;
```

- Each element of the array is a scalar, which is why an individual element of `@fred` is accessed using `$fred[expr]` not `@fred[expr]`. This is admittedly confusing!
- Assigning to an element beyond the current end of the array causes the array to be silently extended. All intervening elements are made the undefined value (looks just like 0).

- A single item may be extracted from an array:

```
$sum += $fred[$i];
```

The index expression will be truncated to an integer before the array is accessed.

- Building an array piece by piece is painful: instead, assign a bracketed comma-separated list of scalar expressions straight into an array:

```
my @fred = ( "hello", 7.1+$a, 17.3, $c );
```

- Inside a list, the `..` operator can be used as in `(1..20)` or `('a'..'z')`.

- If you have a list of single words, for example:

```
my @fred = ( "hello", "there", "how", "are", "you" );
```

- Perl provides the following syntactic sugar:

```
my @fred = qw(hello there how are you);
```

- You can also break up an array into a list:

```
my( $a, $b, $c ) = @fred;
```

- This copies `$fred[0]` to `$a`, `$fred[1]` to `$b` and `$fred[2]` to `$c`. Any remaining elements in the array are ignored. If `@fred` has (say) only 2 elements then `$c` is set to the undefined value.
- An array can be used to soak up the remainder:

```
my( $a, $b, @c ) = @fred;
```

- This gives you a very easy swap operation:

```
( $x, $y ) = ( $y, $x );
```

which takes `y` and `x`, forms them into a two-element list, and assigns the first two elements of that list back into `x` and `y`.

- In summary, Perl arrays act as dynamic arrays, tuples, lists, stacks and queues.

- Some operators behave differently when placed in a **scalar context** or in an **array context**. An array context is one where an array is expected rather than a scalar.

- `<>` is one such operator:

- In a scalar context, eg `$line = <$in>`, it reads a single line.
- In an array context, eg `@line = <$in>`, it reads the *rest of the input*, returning an array of lines - still with all the newlines.
- Fortunately, `chomp @line` chops the newline from every line.

- Similarly, array assignment:

- Assigning an array to another array, eg. `@x = @y`, copies the entire array.
- Assigning an array to a scalar, eg `$count = @y`, means *set \$count to the number of elements in @y*.

- You can force a scalar context when you're not sure what Perl would do by wrapping an expression in the function `scalar()`.

- Declare a hash variable by `my %fred`, such a hash occupies a different namespace from `$fred` and `@fred`.
- A hash stores (*key, value*) pairs - for each string scalar (the *key*), it stores an arbitrary scalar (the *value*).
- A hash literal can be written as a list of pairs with the `key => value` syntax:

```
my %roomno = (
    "dcw" => "225",
    "sza" => "226",
    "iwm" => "228"
);
```

- You can think of a hash as a two-column database table (but stored in memory), indexed on the key:

| Key | Value |
|------|-------|
| dcw | 225 |
| sza | 226 |
| iwm | 228 |
| | |

- Hashes have a highly efficient indexing system (a hash table, hence the name), so you can look up a key's value very quickly.

- To add a single (*key, value*) pair into a hash, do:

```
$roomno{"susan"} = "569";
```

- As a convenience, Perl allows you to omit the quotes on a literal key string, and write:

```
$roomno{susan} = "569";
```

- To check whether a key is present in the hash, use `exists`, eg:

```
print "elvis has left the building\n" unless exists $roomno{elvis};
```

- To delete a single (*key, value*) pair from a hash:

```
delete $roomno{dcw};
```

- Alternatively, the entire hash may be cleared by:

```
%roomno = ();
```

- Our original hash literal example could be written as:

```
my %roomno = ();
$roomno{dcw} = "225";
$roomno{sza} = "226";
$roomno{iwm} = "228";
```

- To retrieve a particular value from a hash, use:

```
my $room = $roomno{$person};
```

If no value has been stored against the key `$person`, the undefined value is returned.

- To process an entire hash, you can use the `keys()` function:

```
foreach my $key (keys(%roomno))
{
    my $value = $roomno{$key};
    print "$key in room $value\n";
}
```

- Note that `keys(%roomno)` builds an array containing the keys of `%roomno`. This could be very large!

- The keys do not come out in any predefined order - certainly not alphabetical order, or insertion order! They come out in an efficient internal order! Because of this, you often see:

```
foreach my $key (sort keys(%roomno))
...

```

- Very occasionally, you only need the **values**:

```
foreach my $value (values(%roomno))
{
    print "room value $value\n";
}
```

- If you need both keys and values, use the `each()` function and a **while** loop to iterate over all the (*key, value*) pairs:

```
while( my($key,$value) = each(%roomno) )
{
    print "$key in room $value\n";
}
```

- Note: `keys()`, `values()` and `each()` produce results in **whatever order Perl likes**.

- Exercise: build a Perl program storing (`username, roomno`) pairs in a hash, which then reads usernames - from stdin, or a text file, as you like - until the end-of-input is reached, and prints the name and corresponding room number for each.

- Then replace the set of names and their associated room numbers with an external data file, reading them in and then building the in-memory hash. (At this early stage in your Perl knowledge, you might need to store usernames and room numbers on adjacent lines in the file).

- Next, replace the hash with a dbm file (with an initialization program that reads the names and room numbers from a text file, and stores them in the dbm-tied hash).

Perl has many special variables (see `perldoc perlvar` for a complete list). Here are a few of the most useful:

- When you invoke one of your Perl programs, you can place *arguments* on the command line, eg:

```
myprog first second third
```

- When you do this, Perl makes the strings `first`, `second` and `third` available in a special array called `@ARGV`. Specifically:

```
$ARGV[0] = "first";
$ARGV[1] = "second";
$ARGV[2] = "third";
```

- As usual, `@ARGV` evaluated in a scalar context gives the number of elements (in the example, 3).

- The array function `shift()` can be used on `@ARGV`:

```
my $arg = shift @ARGV;
```

This sets `$arg` to element 0 of the array, and removes that element from the array, shifting the other elements down one.

- It's up to the program to decide what the strings *mean!*
- If they are filenames to be opened and processed, the *open and process every line in every file* idiom is often used:

```
foreach my $arg (@ARGV)
{
    my $in = new IO::File( $arg ) || next;
    while( my $line = <$in> )
    {
        chomp $line;
        # now process $line
    }
    $in->close;
}
```

- The above pattern (processing several files, not caring where one ends and the next begins), is so common that Perl has a special shorthand:

```
while( my $line = <> )
{
    chomp $line;
    # now process $line
}
```

- Exercise: generalise one of the earlier STDIN or single-file processing programs to take one or more command line arguments using either of these idioms.

- You may find a puzzling shorthand, as in eg2:

```
while( <> )
{
    chomp;
    print "found '$_'\n" if /dun[ck]/i;
}
```

- Where are we storing the line we read?
- What are we chomping?
- What are we matching /dun[ck]/ against (in a case insensitive way)?
- What's that \$_ interpolated into the print?
- \$_ is the *implicit variable*: the *default argument* to many functions:
 - The default variable where <> stores its input line.
 - The default variable that chop and chomp modify.
 - The default variable to match a regex against.
 - The default value to print if none is given.
 - The default foreach variable, as in foreach (@ARGV).
 - .. and many more cases.

- Unix (like most other platforms) has *environment variables*: arbitrary (name, value) pairs, created by setenv NAME value commands in the shell (by convention, environment variable names are upper case).
- To see the current set of environment variables, type env at the command line. A list of NAME=value pairs fly past.
- Once set, environment variables are passed around automatically to every Unix process in the current session. Perl makes these variables accessible via a single hash called %ENV.
- For example, an important environment variable is HOME (the pathname of your home directory). Get this by:

```
my $home = $ENV{HOME} || die "no home?\n";
```

- Other platforms – such as Windows – also have environment variables, Perl on those platforms can access environment variables in the same way, but of course what environment variables exist and what they mean) are different.

- We saw in the first session that we could write:

```
if( $name =~ /^Dun[ck]/ )
```

- This is an example of matching a string against a *regular expression* (or *regex*), as in the Unix filters **sed**, **grep** and **awk**.
- A regular expression is a way of describing a class of similar strings in a very compact pattern notation. In the above example, the match will succeed if the current value of \$name starts with:
 - A capital 'D' [must be the very first character],
 - The lower case letters 'u' and 'n' as the next two characters,
 - then *either* a lower case 'c' or a 'k'.
- A whole regex is (usually) placed inside a pair of '/' signs. Within the slashes, characters are interpreted pretty much like in a double-quoted string. In particular, **variables are interpolated** before pattern-matching occurs.
- A regex is made up of *single character patterns*, *grouping patterns*, *alternation patterns*, *anchoring patterns* and *bracketing patterns*. We'll look at each in turn.

- ‘.’ matches any single character.
- A single printable character matches itself (except meta-characters like ‘.’, ‘*’ etc, which may be preceded by a backslash when you really want to match the character itself!).
- [set] matches any single character in the set. For example, [aeiou] matches any single lower-case vowel.
- Also, the set may contain items of the form a-f, which is a shorthand for abcdef. For example, [a-z#%] matches any single lower-case letter, a hash-mark, or a percent sign.
- If a set starts with a ‘^’ character (eg. [^a-z#%]), the set is negated - the pattern matches any character NOT in the set.
- Several useful character classes are predefined:

| | | |
|----------------|----|------------------|
| Digit | \d | [0-9] |
| Non-digit | \D | [^0-9] |
| Word | \w | [a-zA-Z0-9_] |
| Non-word | \W | [^a-zA-Z0-9_] |
| Whitespace | \s | space or tab |
| Non-whitespace | \S | not space or tab |

- **Sequence of single-character patterns:** matches a corresponding sequence of characters. eg. /[a-z]bc/ matches any lower case letter, followed immediately by a ‘b’, followed immediately by a ‘c’, anywhere in the string.
- **Optional:** ‘?’ makes the previous pattern optional - i.e. match zero or one times. eg. /he?llo/ matches ‘hello’ or ‘hlllo’.
- **Zero-or-more:** ‘*’ makes the previous pattern apply any number of times (from 0 upwards). eg. /he*llo/ matches ‘hlllo’, ‘hello’, ‘heello’ etc. It consumes the maximum number of ‘e’s possible (it’s **greedy**).
- **One-or-more:** ‘+’ means match 1 or more times. eg. /he+llo/ matches ‘hello’, ‘heello’, ‘heello’ etc but not ‘hlllo’.
- If the greediness of ‘*’ and ‘+’ is ever a problem, use *? or +? to consume as few characters as possible.
- A regex can contain several of these operators: eg: /h[uea]*l+o/ matches ‘hlo’, ‘hullo’, ‘hullllllo’, ‘heello’, ‘heuaueaaeuellllllllo’ etc.

- Placing ‘^’ at the start of a regex matches the *start* of the string. Similarly, ‘\$’ at the end of a regex matches the *end* of the string.
- ‘\b’ constrains the regex to match only at a word boundary.
- Without any anchoring, the regex can match anywhere.

There are two main ways of using regexes:

- To check whether a **string matches a regex**. We specify the string to match against using the =~ operator, or the *not match* operator !~:

```
print "<str> matches\n" if $str =~ /h[eua]*l+o/;
```

If a regex match is followed by i, as in /h[eua]*l+o/i, the matching is done case insensitively.

- Secondly, a regex can be used to **search and replace** all occurrences of a regex within a string (again, we specify the string to modify using the =~ operator):

```
$str =~ s/[aeiou]+/a/g;
```

The trailing g makes Perl replace ALL vowel sequences in \$str with ‘a’. Without the g Perl would only replace the first match.

- As a general way of testing regular expressions, I recommend a program like **eg3**:

```
#!/usr/bin/perl
#
# eg3: regex test harness..
#
print "Please enter a string: ";
my $str = <STDIN>;
chomp $str;
print "\nat start : <str>\n";

# test search and replace:
$str =~ s/^\s+//;
print "\nafter s///: <str>\n";

# test pattern match:
print "\n<str> matches hello regex\n" if $str =~ /h[eua]*l+o/;
```

- This whole program exists in order to let you test search and replace and/or pattern matches using a string entered at the keyboard. By the way, s/^\s+// is a useful regex - worth committing to memory - that removes any leading whitespace. Similarly, the regex s/\s+\$/ removes trailing whitespace.
- I strongly recommend that you use this program to test lots of different regexes and their behaviour against various strings.

- A regex of the form `/h[euɑ]*llo|wo+tcha/` matches *either* `/h[euɑ]*llo/` or `/wo+tcha/`. Note that `/a|b|c|g/` should be written as `/[abcg]/` instead for efficiency.
- Brackets may be placed around any complete sub-pattern, as a way of enforcing a desired precedence. For example, in `/so+ng|bla+ckbird/` obviously `bird` is only part of the second alternative (`bla+ckbird`).
- If you meant `/so+ng|bla+ck/` followed by `/bird/`, then write that as `/(so+ng|bla+ck)bird/`.
- If you want a repetition ('+', '*' or '?') of *anything* longer than a single character pattern, you need brackets, as in `/(hello)*/` - `/hello*/` means `/hell/` followed by `/o*/`!
- Brackets have another useful side effect: they tell Perl's regex engine to *remember* what text fragment matched the inner pattern for later reporting or reuse.

- For example, in the code:

```
my $str = "I'm a melodious little soooongbird, hear me sing";
print "found <$1>\n" if $str =~ /(so+ng|bla+ck)bird/;
```

After the match succeeds, the capture buffer variable `$1` contains `soooong` - the part of `$str` matching the bracketed regex.

- Use up to nine bracketed sub-patterns in a single pattern match - capture buffer variables `$1` to `$9` - available for use as soon as the pattern match has succeeded.
- Capture buffers can be used in a search and replace operation:


```
$str =~ s/^(w+)\s+(\w+)/$2 $1/;
```

 which - if there are two or more space separated words at the front of `$str` - swaps the first two.
- Another example: `/first(.*)second/` matches exactly the same strings as `/first.*second/`, but remembers the particular sequence of characters found between `first` and `second` as `$1`.
- If the string contains several occurrences of `first` and `second`, greediness maximises `./` so the regex matches the *leftmost first* and the *rightmost second*:


```
....first...first...second...first...second.....
      ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
      ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
```

- We can also reuse a capture buffer (under the syntax `\1`) to enforce the *same* literal text is found twice in a pattern match:


```
/first(.*)second\1/
```
- This will only match strings like:


```
.....firstXYZsecondXYZ.....
```

 but not strings like:


```
.....firstABCsecondXYZ.....
```
- Test **eg3** out with a variety of inputs and regexes and check you understand how they work.
- If your pattern contains lots of `'/'` characters - while you can write each as `'\/'` - it's easier to change the regex quote character:


```
$str =~ m%"/([~/]*)%/;
$str =~ s!/[~/]!%$!;
```
- Here, the character immediately following `'m'` (for match) or `'s'` (for search and replace) is used as the regex quote character.
- That's a basic overview of Perl regexes; there are loads more features (a few more bizarre ones seem to get added every year or so). See `perldoc perlre` for more details.

```
$str =~ tr/firstcharlist/secondcharlist/[cnds]
```

- `tr` is the character transliterator. It works very like the Unix filter `tr` - turning each occurrence of a character from the first character list into the corresponding character from the second character list.
- eg: `tr/aB/Ab/` uppercases every 'a' and lowercases every 'B'.
- `tr` is rather like a series of regexes that only use character classes - the above example is equivalent to `s/a/A/g` followed by `s/B/b/g`. But `tr` is much more efficient.
- `tr//` is bound to a variable using the `=~` syntax (like regexes).
- Like `s///`, `tr//` also returns a scalar value - a count of how many characters were modified/deleted.
- Let's give some examples:

| | |
|---|---|
| <pre>\$str =~ tr/A-Z/a-z/ \$str =~ tr/xyz/ZYX/ \$str =~ tr/A-Z//d \$str =~ tr/A-Z//cd \$str =~ tr/aeiou/V/ \$str =~ tr/aeiou/V/s \$count = (\$str =~ tr/a-z/a-z/)</pre> | <pre>lowercase every character in \$str. turn every occurrence of x into Z, y into Y and z into X. delete all upper case letters. delete all characters except upper case letters. replace any lower case vowel with a 'V'. replace each sequence of vowels with a single 'V'. Set \$count to the number of lower case letters found in \$str (without changing \$str).</pre> |
|---|---|

- Implement a simplified Perl version of the Unix grep command - take two or more command line arguments (die with a nice usage message if not enough!):
 - First argument: a literal string to search for (use `shift @ARGV` to extract and remove it).
 - All remaining arguments: filenames to search for the above string.
- Use the process-all-lines-in-all-files idiom, and print out the current filename and line if the line contains the search string.
- Add line number counting (reset it for each file) and print out the current line number on matching lines (as well as the filename and line itself).
- Now add in a sanity check to ensure that the search string does not contain any regex meta-characters. Die with a nice helpful message if it does (or escape them first!).

- Prepare an input file containing a list of words, in no particular order, one per line. Write a program to open such a file - take the filename on the command line - read each line, delete leading and trailing whitespace from each line, delete leading or trailing punctuation too, and then print each line (word) out.
- Now make this word-splitter program count word frequencies - do `$freq{$word}++` for each word `$word` you find. After processing all lines, print out a sorted list of frequencies of all the words found - using magic syntax:


```
foreach my $word (sort {$a <=> $b} (keys(%freq)))
```
- Now, wrap the complete make and print frequency table logic in a loop that processes each of the files named in `@ARGV` - emptying the frequency hash for each file.
- Now, modify the program so that `$ARGV[0]` contains a word to search for, and all the rest of `@ARGV` contains filenames to look in.
- For each file, first use your frequency building code to build the frequency table for that file.

- Then print out the filename if the particular word was present in the table (i.e. if the frequency of that word is more than 0!).
- This is now a primitive indexer - you give it a word and a list of filenames (all in one-word-per-line format) and it searches for the word in all the filenames.
- It would be nice to generalise this to normal text files with multiple words per line, but you don't yet know how to split a string apart into pieces. Delay this until next time!
- Now store the frequency arrays to disk, so that next time we could just use the frequency table not have to recalculate it!
- To do this, you'd need to use a Unix DBM file for each file's frequency array. It would be sensible to store the DBM files in a separate directory, to avoid cluttering up the normal directory.
- You'd also need two programs - one to index (or reindex) a list of files, and another to perform a search for a word...