# An Introduction to Perl

Jason Stajich
UC Riverside
@hyphaltip @stajichlab

# Hello World
## A basic script

```
1 #!/usr/bin/perl
2 print "Hello World.\n";
```

; to end the line
**print** to display output (by default will go to your terminal window)
the \n character is for printing a return at the end of what is printed

To run this, create this text file file called 'hello.pl' with this info. Use a text editors such as 'nano', 'emacs', 'vi', then run on the command line:

```
$ perl hello.pl
```

You can make it executable by doing:

```
$ chmod u+x hello.pl
$ ./hello.pl
```

# Primitives

Comments start with #

```
1 print "some code here, followed by comment"; # here is a comment
```

Lines end with a ';'
Variables are used to store and access data. Scalar variables start with a '$' and can contain one thing.
This can be a string, a number, or a reference. (We'll talk about References later)

```
1 $var1 = "apple";
2 $var2 = "orange";
3 $var3 = 10;
4 $var4 = 17.4312;
```

Print can also operate on a list of things.

```
1 print "Hello there ", $var2, "\n";
```

# Strings

Strings are collections of characters. They can be enclosed in single or double quotes. In Perl, many things can be converted to strings too.

```
1 print "You're turning into penguin.\n";
2 print 'You're turning into penguin.\n'; # this will have an error
3 print 'You\'re turning into penguin.\n'; # this will not have an error, but won't have a r
```

- ' Single quotes are literally what is in the quotes, no special characters or variable interpolation is done

- " Double quotes will allow for special characters and interpolation

- ` (backquote) is used to execute programs outside of Perl (more on this later)

# String operators

Strings can be manipulated and compared in several ways.

- The `.` operator concatenates two strings

- `substr` to get a substring , the arguments are the starting string, the offset, the length of the subsequence (if omitted will go to the end), and (optional) the sequence to replace this substring with.

- `reverse` to reverse the order of a string

- `length` returns the length of a string

- `index` searching from left to right, find the first position of a substring in a longer string. `rindex` does the same searching right to left. An optional 3rd argument tells where to start searching from (instead of default of 0) letting you skip ahead.

# String operator: Code

```
 1 $left = "ABCD";
 2 $right = "XYZ";
 3 $complete = $left.$right;
 4 print "$left $right => $complete\n"; # --> "ABCD XYZ => ABCDXYZ"
 5
 6 $str = "AACTTTGGA";
 7 print substr($str,3), "\n";    # --> TTTGGA
 8 print substr($str,6,3), "\n"; # --> GGA
 9 substr($str,6,3,'NNN'), "\n"; # --> replace GGA with 'NNN'
10 print "updated string is $str\n";
11
12 $rev = reverse $str; # reverse the order
13 print $str,"\n",$rev,"\n";
14
15 $len = length $str; # get the sequence length
16 print "DNA is $len bp long.\n";
17
18 $index = index($str,"TTG"); # find substring inside a string
19 print "first TTG is at $index ", substr($str,$index,3),"\n";
```

# Other string operations

- `lc` and `uc` to upper case and lower case a string

- `chop` will remove the last character from a string

- `chomp` will remove the last character from a string IF it is a whitespace

# WHOA! Too many new things!

How can I get help about these functions? Perldoc is your guide. You can find perldoc in your terminal with the following commands

```
1 perldoc -f substr
2 perldoc -f length
3 perldoc -f reverse
4 perldoc -f index
```

Here are the links to online Perldoc.

- substr

- length

- reverse

- index

- rindex

# Interpolation

<u>Interpolation</u> is when a variable's value is evaluated and substituted, for example when you want to insert a value into a string.

```
1 $fruit = 'apple';
2 print "my favorite fruit is $fruit.\n";
3 $fruit = 'grape';
4 print "my favorite fruit is $fruit.\n";
5
6 $fruit = "$fruit-$fruit";
7 print "my favorite fruit is $fruit.\n";
```

This will produce

```
my favorite fruit is apple.
my favorite fruit is grape.
my favorite fruit is grape-grape.
```

# Naming variables

Variable names must contain alphanumeric characters. You can name variables how you like, note that variables cannot start with a number, but can contain numbers.

```
1 $xyz = "ATGCAGTGA"; # not very descriptive
2 $protein = "ATGCAGTGA"; # misleading
3 $dna_sequence_variable = "ATGCAGTGA"; # needlessly long
4 $sequence = "ATGCAGTGA"; # better
5 $dna = "ATGCAGTAGA"; # even better
```

- Often we name variables by their use or to describe their use.

- Other convention `$i,$j,$k` are often counters.

- `$a` and `$b` are special variable names in some contexts in Perl so try to avoid using these.

- Some conventions are to use underscores_to_separate_spaces or CamelCase

# Numerics

Numbers can be integers, floating point, scientific notation. They can be initialized and computed.
Strings can also be converted to numbers.

```
 1 $x = 10;
 2 $y = 2.2;
 3 $z = 1/3;
 4 print "$x is x, z is $z\n";
 5 $sum = $x + $y;
 6 print "sum = $sum\n";
 7
 8 $string = "12";
 9 $sum = $string + $x;
10 print "sum = $sum\n";
11
12 $sum += 5; # can add to an existing in one line
13 $sum /= 5; # or divide by
14 print "sum = $sum\n";
```

This will produce

```
1 10 is x, z is 0.333333333333333
2 sum = 12.2
3 sum = 22
4 sum is 5.4
```

# Assignment of values

Values from one variable can be copied to another.

```
 1 $dollars = 50;
 2 $cents   = $dollars;
 3 $cents   *= 100;
 4 print "$dollars dollars = $cents cents\n";
 5
 6 $x = 100;
 7 $y = $x;
 8 print "x=$x y=$y\n";
 9 $x = 50;
10 print "x=$x y=$y\n"; # won't change the value of y to update x
```

Will produce

```
1 50 dollars = 5000 cents
2 x=100 y=100
3 x=50 y=100
```

# Playing it Safe

So far have shown very basic code. Without any warnings turned out. Only syntax or system errors will cause the program to stop without extra options turned on. By default Perl will not warn you about empty or undeclared variables. To be a better programmer you want to use the following best practices to write better Perl code.

1. use strict

2. use warnings

3. use `my` to declare variables. This declares them in a particular scope.

# Some examples of errors

```
 1 $x = 7;
 2 print "$x and $y\n"
 3
 4 use warnings;
 5 $x = 7;
 6 print "$x and $y\n"
 7
 8 use strict;
 9 use warnings;
10 $x = 7;
11 print "$x and $y\n";
12
13 use strict;
14 use warnings;
15 my $x = 7;
16 my $y;
17 print "$x and $y\n";
18
19 use strict;
20 use warnings;
21 my $x = 7;
22 my $y = 'PRP8';
23 print "$x and $y\n";
```

# At the top

Typical boilerplate for your Perl scripts then should usually look like this:

```
1 #!/usr/bin/perl
2 use strict;
3 use warnings;
```

# Some typical errors

Uninitialized variable, $y had no value.

```
1 use warnings;
2 use strict;
3 my ($x,$y);
4 $x = $y+10;
```

Syntax error – missing semicolon on line 3.

```
1 use warnings;
2 use strict;
3 my ($x,$y)
4 $x = $y+10;
```

Undeclared variable, forgo to declare $y before we used it.

```
1 use warnings;
2 use strict;
3 my ($x);
4 $x = $y+10;
```

# Other numeric operators

- `+,-,*,/` for add, subtract, multiply, and divide

- `=` is the assignment variable to assign a value, not be confused with `==` which is for testing if two values are equivalent

- Assignment and an arithmetic operation for updating a variable, so `$val += 10` adds 10 to the current value

- The modulus operator `%` does a division and returns the remainder. So `9 % 2` is 1 because 2 will divide into 9, 4 times, leaving a remainder of 1.

- Exponent with **, so 2^5 is `2**5`

# Printing Output

Sending information back out from the program is important. For all the computational you will be doing, you want to know the answer.

- `print` will send output to STDOUT by default. You can explictly choose the stream by providing it as in `print STDOUT "hello world\n";` or `print STDERR "hello world\n"; Other streams can be defined such as an output file, we'll cover that in lecture 2.

- `printf` will allow for formatted printing using a variety of options, it requires multiple arguments, one is a formatting string, the rest are the specific data to be interpolated into the formatting string.

- `sprintf` is like printf except it returns a formated string rather than printing it to a stream.

    - %d – for integers

    - %f – for floating point, %.2f – to specify number of significant figures, %5.2f to format the width of the string (padding with spaces).

    - %s – for strings

    - %g – for scientific notation

# printf examples

```
1 $str = "MAYWRCILR";
2 printf "This '%s' string is %d letter long\n", $str, length($str);
3 printf "This is number %d\n", 17;
4 printf "This is a floating point number %f\n", 4.25;
5 printf "This is a floating point number as an integer %d\n", 4.25;
6 printf "This is a floating point number with three signif digits %.3f\n", 12/11;
7 printf "This is scientific notation %g\n", 21213333;
8 printf "This is scientific notation with 2 signif figs %.2g\n", 21213333;
9 printf "Formatted length '%4d' and '%7.2f'\n", 21, 2/7;
```

Produces these results

```
1 This 'MAYWRCILR' string is 9 letter long
2 This is number 17
3 This is a floating point number 4.250000
4 This is a floating point number as an integer 4
5 This is a floating point number with three signif digits 1.091
6 This is scientific notation 2.12133e+07
7 This is scientific notation with 2 signif fig 2.1e+07
8 Formatted length '  21' and '   0.29'
```

# Danger Will Robinson!

- Warnings can be useful to print when something is unexpected.

- One can use the `print STDERR` to print to the error stream.

```
1 print STDERR "There was a problem in this program"
```

I also find that using the `warn` command is much more useful.

```
1 if( $error_condition == 1 ) {
2     warn("There was a problem sir!\n");
3 }
```

# Kill them all

Sometimes you want to exit your program. There may be an unresolvable error, or you want to just finish right there. Two commands are useful for this.

- `exit` which will exit the program. It takes an optional numeric to return to the operating system, but usually you just use it alone.

- `die` which will print a warning message and exit, or if you fail to include a warning message it will report the line of the error, which can be helpful in debugging if you have lots of code.

# Using die

Here is a script snippet from the Unix & Perl book to demonstrate testing if a sequence could encode a protein (by being divizable perfectly by 3).

```
use strict;
use warnings;
my $dna = "ATGCACGTAGGTAACACTGACTGAA";
my $length = length($dna);
die "$length is not a valid length\n" if (($length % 3) != 0);
```

# Other special characters

Some special string characters

- `\t` for tab

- `\n` for newlines, on windows this is `\r` and on mac this is `\r\n`. But please just use `\n`

- `chr` will convert an ASCII number, e.g. `print chr(97),"\n"` will print `a`.

Perl has several special variables `$_` – *the implicit variable, default variable for matching* `$,` – is the default character to separate values when print an array as a string

# Arrays and Lists

In addition to the scalar variables (`$var`) that can store numbers or strings, there are also variables for multiple items in lists. The array variables start with `@` and represent multiple items.

```
1 @gems = ('diamond','ruby', 'quartz');
2 print "@gems\n";
3 $one = 'apple';
4 $two = 'pear';
5 @fruit = ($one,$two);
6 print "@fruit\n";
```

Arrays can be indexed by a number, so the 1st item in the array can be obtained with `$array[0]`, the second with `$array[1]` and so on. You can ask for the last item in an array with `$array[-1]`, the second to last with `$array[-2]` and so on.

Perl will allocate new memory for you on the fly as your array gets bigger – so if you ask for `$array[100]` it will create 0–99.
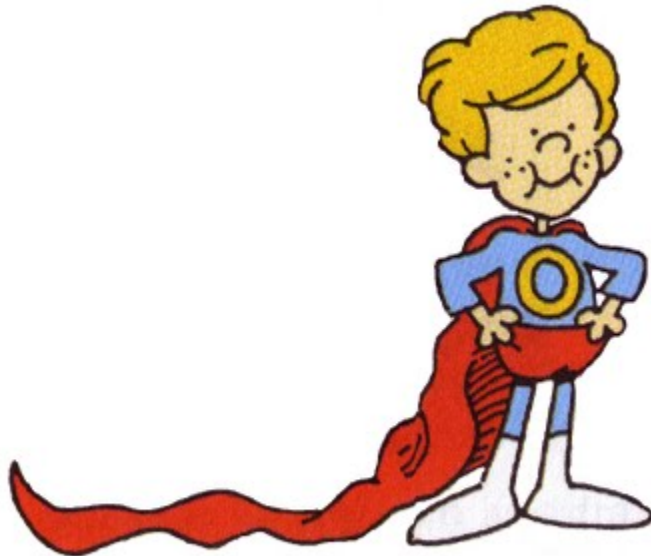
If you want to crash perl you can do something as simple as (depending on how much memory your machine has)

```
1 $x[10000000000000] = 1;
```

# My hero, Zero

In computer science, when starting to count, the first value is always 0 not 1. So lists and strings are indexed starting with 0. So if you want the 1st item in a list, you need to ask for the `[0]` one. Similarly, to get the second character in a string you ask for the one at position `1`. So `substr($string,1,1)` will return the 2nd character in a string.

```
['a','b','c','d']
  0   1   2   3
```

# Array operations

- Length of an array is obtained by treating it like a scalar variable, either `$len = @gems;` or more typically `$len = scalar @gems;`

- `sort @array` will order a list alphanumerically. It can be ordered numerically with some additional code, we'll discuss later

- `join($joincharacter,@array)` transform an array into a string, joining together with the intervening character

- `@array = split("\t",$string)` transform a string into an array, splitting on a particular character

- Modify and array adding or removing items

  - `$last = pop @array` – remove the last item from the array and return it

  - `$first = shift @array` – remove the first item from the array

  - `push @array, $enditem` – add an item to the end of the list

  - `unshift @array, $firstitem` – add an item to the beginning of the list

- `splice` – a cool operator to get out something from the middle

# Array operations: Code

```
 1 my @array = ();
 2 push @array, 'yellow';
 3 unshift @array, 'red';
 4 print join(" ", @array),"\n";
 5 push @array, 'green', 'purple','blue';
 6 print join(" ", @array),"\n";
 7 print $array[2], "\n"; # what do you think will be printed here?
 8
 9 my $str = "AACA-AG-TTTG-TACA";
10 my @bases = split('-',$str);
11 print "bases are @bases\n";
12 $str = "AACA--AG-TTTG-TACA";
13 my @bases = split('-',$str); # this doesn't quite work if there are adjacent gaps
14 print "bases are @bases\n";
15
16 $str = "AACA--AG-TTTG-TACA";
17 my @bases = split(/\-+/,$str); # this way will work using patterns
18 print "bases are @bases\n";
19
20 @array = (1..10);
21 print "3 elements starting with 6th one are ", join(" ",splice(@array, 5,3)),"\n";
22 splice(@array,2,2,'100');
23 print join(" ", @array), "\n";
```

# Sorting

A useful operation on a list is the ability to sort the data. The default sort mechanisms.

This is the place where some special variables `$a` and `$b` are used. They have special meaning as a sort is done by comparing pairs of numbers. So you just have to define an operation which can determine for two values which is larger or smaller (or if they are equal).

The `<=>` operator is great for this and operates on numbers. It returns −1 if the left is smaller, 1 if the left is larger, and 0 if the two values are equal. The equivalent operator for alphanumeric is called `cmp`. You can see more on operators at `perldoc perlop`.

```
1 @nums = (20,1,18,3.2,71,53);
2 print sort @nums; # print alphanumerically
3 print sort { $a <=> $b } @nums; #print smallest to largest
4 print sort { $b <=> $a } @nums; #print largets to smallest
5 print sort { length($a) <=> length($b) } @nums; #print by num of digits
```

# Convience method for array init

The `qw` (quote words) operator is useful for initializing a list, it separates items by white space from a list and converts into an array simplifying the code.

```
1 @wintermonths = qw(Dec Jan Feb);
2 # instead of
3 @summermonths = ('Jun','Jul','Aug');
```

# Hashes or Associative Arrays

Hashes are like arrays, but they are indexed by strings instead of numbers. Hashes use '{}' to index instead of '[]' in arrays. Variables start with `%`.

They are like dictionaries which means that when you want to lookup something, you look it by the word, not the order.

Arrays are ordered lists -- you can get the 5-th item of the list. However, for a hash, there is no order, so you want to lookup an item by the key.

# Hash operations: Code

```perl
 1 my %months2season;
 2 $months2season{'Jan'} = 'Winter';
 3 $months2season{'Aug'} = 'Summer';
 4 $months2season{'Sep'} = 'Fall';
 5
 6 my $month = 'Aug';
 7 print "$month is in $months2season{$month}\n";
 8
 9 my %bball_teams = ('Los Angeles' => 'Lakers',
10                    'Phoenix' => 'Suns',
11                    'Charlotte' => 'Bobcats');
12
13 print join(",",sort keys %bball_teams), "\n";
14 print join(",", values %bball_teams), "\n";
```

# Hash operations

- `keys %hash` will return a list of all the keys in the hash

- `values %hash` will return a list of all the values in the hash

- `each %hash` is used to return all the key/value pairs, however this really only useful when we get to loops (next lecture)