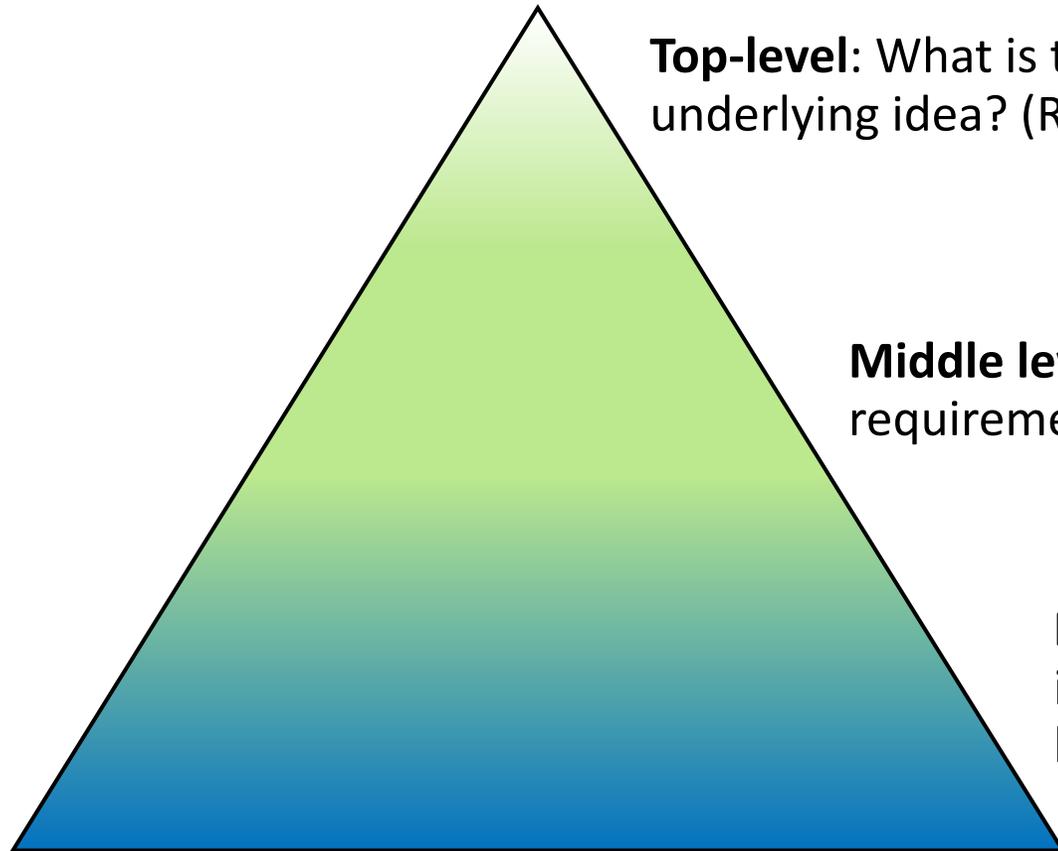




Bridging the gaps in bioinformatics/Getting started in bioinformatics

Why and how to document your code

Documentation can be done at multiple levels



Top-level: What is this code for? What is the underlying idea? (README/Wiki)

Middle level: How can I use it? (Usage info, installation requirements, scope, options)

Deep level: How does the code work? What is the purpose of each chunk? (Comment lines and structure)

Why should you care about documentation?



Transparency

- Colleagues using the code will be able to point out issues, if they understand what the code is doing
- Colleagues with other expertise (wetlab, epidemiology etc) may also be able to point out issues you are not aware of

Why should you care about documentation?



Transparency

- Colleagues using the code will be able to point out issues, if they understand what the code is doing
- Colleagues with other expertise (wetlab, epidemiology etc) may also be able to point out issues you are not aware of

Collaboration

- Give your colleagues a chance to contribute and improve on the code
- Give your colleagues a chance to take over from you

Why should you care about documentation?



Transparency

- Colleagues using the code will be able to point out issues, if they understand what the code is doing
- Colleagues with other expertise (wetlab, epidemiology etc) may also be able to point out issues you are not aware of

Collaboration

- Give your colleagues a chance to contribute and improve on the code
- Give your colleagues a chance to take over from you

Reproducibility

- This is how we do science!

Why should you care about documentation?

Transparency

- Colleagues using the code will be able to point out issues, if they understand what the code is doing
- Colleagues with other expertise (wetlab, epidemiology etc) may also be able to point out issues you are not aware of

Collaboration

- Give your colleagues a chance to contribute and improve on the code
- Give your colleagues a chance to take over from you

Reproducibility

- This is how we do science!

Helping yourself out

- It is surprisingly easy to forget the details of code you wrote a while ago..

README!

- Format of README file is not fixed. Content depends on your project
- Start with overview, then go more detailed further down.
- The worst README is the one that doesn't exist!

Update your README file continuously while you develop your code!

Usage information

Well-written code has integrated usage information.

- Something useful should happen if you execute the software without arguments (or with the -h or --help flag).

```
$ bcftools

Program: bcftools (Tools for variant calling and manipulating VCFs and BCFs)
License: GNU GPLv3+, due to use of the GNU Scientific Library
Version: 1.14 (using htlib 1.14)

Usage:  bcftools [--version|--version-only] [--help] <command> <argument>

Commands:

-- Indexing
   index          index VCF/BCF files

-- VCF/BCF manipulation
   annotate       annotate and edit VCF/BCF files
   concat         concatenate VCF/BCF files from the same set of samples
   convert        convert VCF/BCF files to different formats and back
   isec           intersections of VCF/BCF files
```

Script level documentation

- Keep a clear and logical structure.
- Name your variables properly
- Use comments (#)

get_ambiguous_positions.py

```
1  #!/usr/bin/env python3
2  import sys
3  import os
4  from Bio import SeqIO
5  from Bio.Seq import Seq
6
7
8  #Usage: python length_fasta.py sequences.fasta
9
10 #List of characters, which will not be counted as ambiguous
11 non_amb = ["A","C","T","G","a","c","t","g","N","n","-"]
12
13 #Print a header of output
14 print("Seq_id\tPosition\tAmb")
15
16 #Read in fasta-file, and loop over sequences (first "for loop")
17 fasta_file = sys.argv[1]
18 for seq_record in SeqIO.parse(fasta_file, "fasta"):
19     #Get sequence id, sequence and sequence length
20     seq_id = seq_record.id
21     seq_str = str(seq_record.seq)
22     seq_length = len(seq_str)
23     #Loop over each character in the sequence, and check if it exists in the "non_amb" list. If not, print sequence-id, position, and character
24     position=1
25     for character in seq_str:
26         if character not in non_amb:
27             print(seq_id, position, character, sep="\t")
28             position = position + 1 #remember to add 1 to the position variable, as you move through the sequence string. Also when the character is not ambiguous!
```

Code evolves

You write a script, and all is well. But then:

- There is a bug
- You need another usage option
- Or an additional script
- There is a dependency that got updated
- Change in wetlab

You write a script, and all is well. But then:

- Turns out that last add-on was a really bad idea. Must be changed. Meanwhile, users should go back to the original version. Did you save it somewhere?

=> You need version control software

Acknowledgements

The creation of this training material was commissioned by ECDC to Statens Serum Institut (SSI) with the direct involvement of Kirsten Ellegaard