# Sequence alignments (1)

Why create them?

How are they stored?

How can they be manipulated?
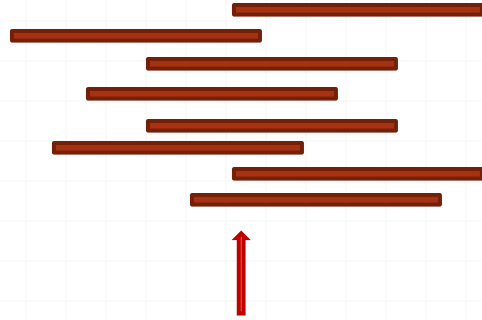
# This presentation

o - is intended to :

    o give a brief overview of types of sequence alignment file formats

    o provide brief background useful for the SAMTOOLs tutorial run in the practical sessions

        o Course notes for tutorial : http://biobits.org/samtools_primer.html

o There will be more detailed presentations / sessions on both the above topics
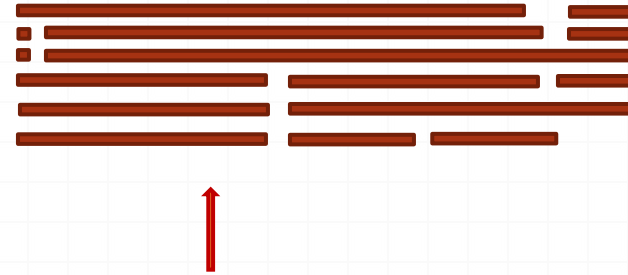
IFR Institute of Food Research

# Sequence alignment data

- What is the purpose of creating and storing a sequence alignment?
- There are multiple purposes
- Different aims
- These aims are associated with different ways of representing the alignment
- And thus different **alignment file formats**

*Bite-sized Bioinformatics. J. Walshaw, GHFS.*

IFR Institute of
Food Research

- Related but different aims, meanings and file formats

- Sequence read alignment ("assembly")

- Multiple protein or nucleotide sequence alignment

o Each nucleotide position (column) represents multiple copies of the same base of an original sequence (e.g. genome sequence)

o Each position (column) represents a **homologous** nucleotide (or amino acid).

o Sequences are evolutionarily related (homologous) sequences, typically from different organisms, and/or multiple members of a gene family

o Gaps represent insertions/ deletions

*Bite-sized Bioinformatics. J. Walshaw, GHFS.*

IFR Institute of Food Research

# "Traditional" alignment of 2 or more sequences

o **Pairwise alignment and multiple alignment**

o From an algorithmic point of view:

o There are some significant differences between methods used to:

   o compare and align only 2 sequences (**pairwise alignment**)

   o compare and align 3 or more sequences (**multiple alignment**)

o Some pairwise alignment tools produce their own particular alignment formats

   o E.g. some programs in the versatile EMBOSS software package

# Example pairwise alignment

## from EMBOSS `needle`

### a **global alignment**, i.e. complete lengths of both sequences are aligned

```
#=======================================
#
# Aligned_sequences: 2
# 1: P0A1L1
# 2: Q0PBF7_CAMJE
# Matrix: EBLOSUM62
# Gap_penalty: 10.0
# Extend_penalty: 0.5
#
# Length: 265
# Identity:      27/265 (10.2%)
# Similarity:    49/265 (18.5%)
# Gaps:         140/265 (52.8%)
# Score: 32.0
#
#
#=======================================

P0A1L1            0 ------------------------------------------------      0

Q0PBF7_CAMJE      1 MRLLVLFFLILPLYSVELISYNIYDRNDRVDLMLSFDNAYNGKISQKKEK     50

P0A1L1            0 ------------------------------------------------      0

Q0PBF7_CAMJE     51 NLTLLTFSDLTYSKDELKELNSQLVDKISISSKNNNTYIMLQNKQNINLE    100

P0A1L1            1 -------------MMKTEATV-SQPTAPAGS--PLM--QVSGALIG---     28
                                 :.:.:|.: |.||..|.:  .||  ..|.:|.|
Q0PBF7_CAMJE    101 LSSINDKFGVRIRAIEQGKANIESAPTTTANNSQELMPKPKSTSLEGYDY    150

P0A1L1           29 --------IIALILAAAWVIKRMGFAPKGNSVRGLKVSASASLGPRERVV     70
                            |:.:::|...|..|:........|..|...:....|....::|
Q0PBF7_CAMJE    151 TNYILVMLILVILLIVLWWFKKTMVYKNNNVSRDFTMIFQRFLDKNNQLV    200

P0A1L1           71 IVEVENARLVLGVTASQINLLHTLPPAENDTEAPVAPPADFQNMMKSLLK    120
                       :.:..|.|..:.:.:.|:.|.....|.|.........:|.:.:..|
Q0PBF7_CAMJE    201 VFDHANKRYTMIIGNSNVLLESIEIPEEQTIKHTEKTEKNFDSFFEENKK    250

P0A1L1          121 RSGRS----------     125
                        |....
Q0PBF7_CAMJE    251 RIQNLIEQRQKGKKS    265
```

IFR Institute of Food Research

# "Traditional" alignment of 2 or more sequences

*O* Typically, DNA sequences representing (complete) genes

*O* Or protein sequences

*O* A principle purpose is to identify the parts of the sequence which are the same and which are different

*O* I.e. identify **conservation** and **divergence**

*O* This can highlight regions of potential functional importance

*O* Can be especially informative if you are dealing with several/many sequences (i.e. a **multiple alignment**)

*Bite-sized Bioinformatics. J. Walshaw, GHTS.*

IFR Institute of Food Research

FASTA

MSF

Stockholm

Some multiple sequence alignment formats – these are all **plain text files** ("flatfiles")

Numerous other formats exist

CLUSTAL format is widely used

26/10/2016

IFR Institute of Food Research

# Viewing multiple alignments

O Example display from a
program (Jalview) which
reads the flatfiles and creates
a formatted/annotated view

IFR Institute of Food Research

- Example BLAST output
  - A list of **pairwise** alignments
  - One sequence (the Query) is always the same
  - In each case, the other sequence is the 'hit' (Subject)
  - N.B. these are **local** alignments (i.e. only of matching segments, not necessarily the whole sequence)
- More compact BLAST output formats are available
- These are **flatfiles** (plain text)
- Again, when interpreted and displayed by other software they may appear very different
- E.g. on a website which provides a BLAST service

```
> tr|A0A031GPC8|A0A031GPC8_9BURK Flagellar biosynthesis protein
FliO OS=Janthinobacterium lividum GN=fliO PE=4 SV=1
Length=196

 Score = 45.8 bits (107),  Expect = 0.009, Method: Composition-based stats.
 Identities = 16/114 (14%), Positives = 40/114 (35%), Gaps = 5/114 (4%)

Query  143   TSLEGYDYTNYILVMLILVILLIVLWWFKKTMVYKNNNVSRDFTMIFQRFLDKNNQLVVF  202
             +            I  ++ ++ LLI L WF K    K    +   ++    L    ++V+
Sbjct  79    PASSAGSLLQTIFALMFVLALLIGLAWFMKRYGPKVMGGNNKMRVVSSLNLGGRERIVLV  138

Query  203   DHANKRYTMIIGNSNV-LLESIEIPEEQTIKHTE----KTEKNFDSFFEENKKR   251
             + A++   +    +  L ++  E   +          NF  + ++  ++
Sbjct  139   EVADQWIVVGASPGRINALATMPRQEGDLPQLATAQNGPAAANFSEWLKQTIEK  192


> tr|A0A0B1REZ0|A0A0B1REZ0_9ENTR Flagellar assembly protein FliO
OS=Pantoea rodasii GN=QU24_01715 PE=4 SV=1
Length=131

 Score = 44.6 bits (104),  Expect = 0.009, Method: Composition-based stats.
 Identities = 15/98 (15%), Positives = 40/98 (41%), Gaps = 2/98 (2%)

Query  156   VMLILVILLIVL-WWFKKTMVYKNNNVSRDFTMIFQRFLDKNNQLVVFDHANKRYTMIIG  214
             V+ ++V+L++   W  K+         ++   +   + +  ++V+ D A+ R  + +
Sbjct  29    VLAVIVLLILACGWLAKRLGFAPKTVNTQALKISASVQVGRQERVVIVDTADARLVLGVT  88

Query  215   NSNVL-LESIEIPEEQTIKHTEKTEKNFDSFFEENKKR  251
             +  L S+     + +      ++F    F+    KR
Sbjct  89    AQQITHLHSLPPVPPEELASNSVAPQDFRQLFQNLVKR  126


> tr|A0A090U6M1|A0A090U6M1_9ENTR Flagellar biosynthesis protein
FliO OS=Citrobacter farmeri GTC 1319 GN=fliO PE=4 SV=1
Length=124

 Score = 44.2 bits (103),  Expect = 0.010, Method: Composition-based stats.
 Identities = 14/84 (17%), Positives = 30/84 (36%), Gaps = 0/84 (0%)

Query  168   WWFKKTMVYKNNNVSRDFTMIFQRFLDKNNQLVVFDHANKRYTMIIGNSNVLLESIEIPE  227
             W  K+    + +R   +    L     ++V+ D  + R  + +  SN+ +       P
Sbjct  37    WVIKRLGFSPKGSHTRGLKVSASTSLGPRERVVIVDVEDARLVLGVTASNISVLHTLPPA  96
```

# "Traditional" alignment of 2 or more sequences

○ Strictly speaking, an alignment is only meaningful if the sequences are homologous (related by descent from a common ancestor)

○ This common descent applies to the individual nucleotide (or amino acid) positions –

  ○ i.e., equivalent bases/ amino acids are lined up

○ This also enables inferences of evolutionary events:

  ○ Substitutions

  ○ Insertions

  ○ Deletions

# "Traditional" alignment of 2 or more sequences

0 But the process of alignment can help to **determine whether the sequences are homologous or not**

0 This principle applies to **sequence similarity search** methods

  0 E.g. **BLAST**

0 - in which a single query sequence is compared, to many sequences in a database, one at a time

0 i.e. **many <u>pairwise alignments</u>** (all involving the query sequence)

0 Each alignment has a score

0 Only those with a sufficiently high score are treated as **<u>hits</u>**

0 The rest are ignored

IFR Institute of Food Research

- One way of very briefly summarising an alignment is to quote a single metric such as:
  - 'percentage identity'
  - 'percentage similarity'
- However, these are not absolutes, and depend on **how** the sequences are aligned
  - I.e., which method and scoring parameters are used
- These, and other aspects of alignment will be described further in a future session

Principles of sequence alignment

*Bite-sized Bioinformatics. J. Walshaw, GHFS.*

IFR Institute of Food Research

# Other reasons for pairwise alignment

o Comparing 2 DNA sequences (or 1 RNA sequence and 1 DNA sequence) –

o where one is expected to be a (often very small) fragment of the other

o E.g.:

o Genome **re**sequencing

    o - and *de novo genome sequencing*

o RNASeq mapping

?

# Genome resequencing

o You already have a **reference genome sequence** of organism X

o The new project is to sequence the genome of organism Y

    o Might be a closely related species

    o Or a different strain of the same species

o Assumption is that the 2 genome sequences are the same or very similar along most of their length

o ▬ Each sequence read from genome sequence Y can be "mapped" to the equivalent position in sequence X

IFR Institute of Food Research

# Genome resequencing

- N.B. depending on the circumstances, it may be just as good (or preferable) to do *de novo* assembly of genome Y
  - And then compare the whole of Y to X
  - This "comparative genomics" can be rather more complex than just a pairwise alignment of X and Y
- Resequencing of many genomes – or particular parts of genomes – helps to identify variations between strains in a species
- And variation among individuals in a population
- E.g. identification of SNPs
  - → ***variant calling***

*Bite-sized Bioinformatics. J. Walshaw, GHFS.*

IFR Institute of Food Research

# RNASeq

0 mRNA transcripts are sampled and sequenced

0 The object is to map each "read" to a reference genome sequence

0 Some reads won't match perfectly, due to sequencing errors

   0 - and even some biological errors

0 This is more complicated in eukaryotes than prokaryotes

   0 Due to the presence of introns (present in reference genome, absent from reads) and splice variants

0 One approach to "community RNASeq" (metatranscriptomics) involves mapping reads to several/many reference genomes

   0 Some reads will originate from genomes which are absent from the reference set

*Bite-sized Bioinformatics. J. Walshaw, GHFS.*

IFR Institute of Food Research

# Storing information on mapped reads - economically



- *Is there any point in explicitly recording all the bases which are the **same** as the reference?*

o Each read can be recorded by specifying:

o Its start position relative to the reference genome sequence

o Its length

o All the differences between the read and the reference

IFR Institute of Food Research

# *De novo* genome sequence assembly



• *Again, it should be necessary to record only positions and differences of each read, relative to the consensus*

o Similar principles

o But no reference genome sequence to compare the reads with

o Various algorithms, which in essence compare the reads with each other

o Produces "contigs"

o A consensus sequence can be produced from each contig

IFR Institute of Food Research

# How can mapping locations and differences be encoded?

0 Various approaches

0 Some have 'evolved' beyond their original purpose

0 **CIGAR** format ("Compact Idiosyncratic Gapped Alignment Report")

   0 (and it *is* idiosyncratic)

0 CIGAR is used in **SAM** format ("Sequence Alignment/Map") files

0 SAM, has a "binary" (compressed) equivalent, BAM

   0 **SAM and BAM are still very widely used**

0 A more recent development is **CRAM** format

   0 Stores only differences compared to reference – **minimal storage space**

   0 CRAM files are **not** plain text; you need utilities (in the CRAMTOOLS software) to view them – displayed in e.g. FASTQ or SAM format

# What do you need to know?

- **So how much do you need to know about these formats?**
  - A little awareness of the concepts will be useful.
  - Do not be concerned with understanding all the details.
    - (Probably not many people do understand them all…
    - the full SAM specification has many details, some more documented than others)
  - There have been many approaches to **compressing** this kind of sequence data, due to the ever-increasing sizes of data sets
    - Compression of SAM/BAM in particular
    - Compression of FASTQ data in general
    - Beyond the scope of this topic

IFR Institute of Food Research

# CIGAR and SAM in brief

A quick overview for now.

There are some more details in another presentation

SAM, pileup and related formats

*Bite-sized Bioinformatics. J. Walshaw, GHFS.*

**IFR** Institute of Food Research

# CIGAR format

- One CIGAR 'string' per alignment (read ↔ reference):
  - Example: '**6=1X7=2X6=**'
- Original CIGAR specification was designed to specify for each alignment:
  - **How long the alignment is**
  - **Where the insertions and deletions (indels) are**
  - And some other attributes like 'clipping', 'padding' (ignore this now)
- Later additions to CIGAR also permitted specification of:
  - ***Where* mismatching bases occur**
  - But ***NOT*** what those mismatches actually are
  - -it was simply not the original purpose of the tools which used CIGAR

IFR Institute of Food Research

# SAM format

0 Each read is represented on one line (1 line = 1 'record')

0 Each line has a field specifying the position on the reference sequence

0 And a field in CIGAR format

0 Then it gets complicated…

0 ***Also additional <u>optional fields</u>*** – these can be used to state:

   0 explicit base differences (the optional `MD` **field**)

   0 (and many other things)

0 SAM format **can** explicitly state the sequence of each read

   0 And often does

   0 But does **not** have to

*Bite-sized Bioinformatics. J. Walshaw, GHFS.*

IFR Institute of Food Research

# That all seems a bit messy....?

## So why is it like this?

IFR Institute of Food Research

# Evolution of bioinformatics formats

0 This is an example of bioinformatics format specifications which have become revised over the years

0 The original design may have been to achieve something very specific

0 Additional functionality has been enabled by additional information included in the data files

0 New fields "bolted on" to a simpler, earlier spec

0 Complete revisions are generally uncommon because of the need to retain backward compatibility

IFR Institute of Food Research

# Some context

0 Example: You want to know **how many reads** are aligned to each part of the genome

*Bite-sized Bioinformatics. J. Walshaw, GHFS.*

IFR Institute of Food Research

# Some context

0 Example: You want to know **how many reads** are aligned to each part of the genome

0 E.g. transcriptomics (RNAseq) – thus, gene expression levels

    0 Which genes are most represented?

    0 You don't care if there are a small number of mismatches due to sequence errors

    0 You don't even care what the sequences of the alignments in each region actually *are*

    0 **All you want to know is where the alignment/mapping program (whatever that may be) mapped the reads**

    0 **Only minimal CIGAR/SAM information is required for this**

    0 (N.B. SAM format can optionally store scores produced by the alignment program)

*Bite-sized Bioinformatics. J. Walshaw, GHFS.*

**IFR** Institute of Food Research

# Summary (1)

o There are many different file formats to store alignment data

o These can be read and created by many different software tools

o Nearly all of these formats are plain text ("flatfile")

    o BAM is not, but is a compressed version of SAM

    o Various other compressed, "binary" formats exist, e.g. CRAM

o These file formats are used for a variety of purposes

    o Which data types are stored in any one file can depend on the purpose

o It is **not** necessary to understand the detailed specifications of these data file formats

o But it is important to know:

    o ***what kind of data* can/cannot be contained in these files**

*Bite-sized Bioinformatics. J. Walshaw, GHFS.*

IFR Institute of Food Research

# Summary (2)

- SAM and BAM data files are principally manipulated with the SAMTOOLS software package
- Various other software can read SAM/BAM files
  - E.g. visual browsers, like the Integrated Genome Viewer
- External link to SAMTOOLS tutorial at BIOBITS:
  - http://biobits.org/samtools_primer.html
- If you are interested in learning more details of SAM, CIGAR etc., these are provided in a further presentation

SAM, pileup and related formats

IFR Institute of Food Research