



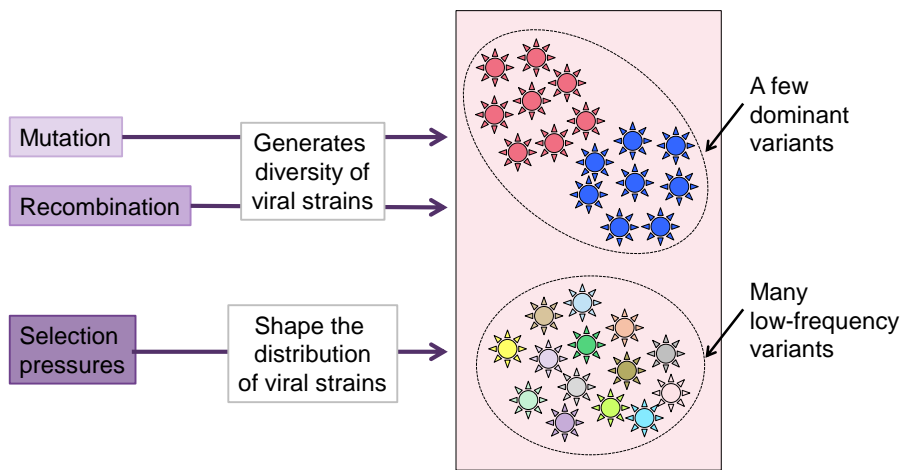


Using sequencing to understand HIV diversity

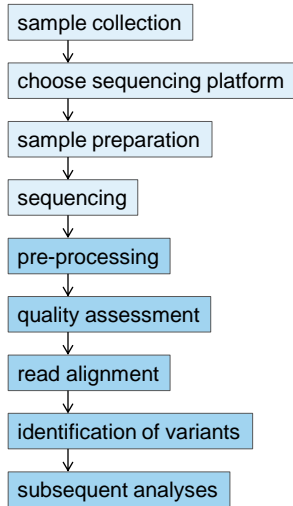
Vanessa Venturi | Infection Analytics Program

Using sequencing to understand HIV diversity  

HIV quasispecies

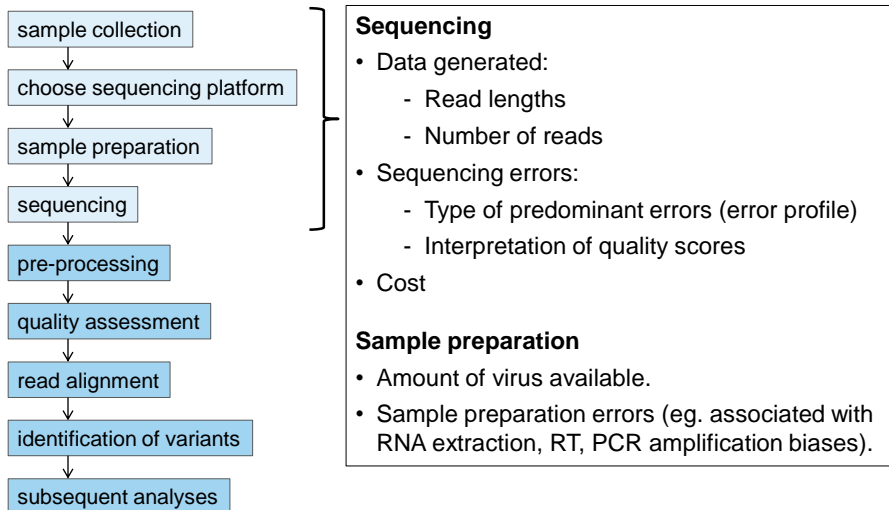


Basic workflow



3

Experimental considerations



4

Comparison of sequencing technologies

	1 st -generation: Conventional	2 nd -generation: Next-generation	3 rd -generation: Single-molecule
Technologies	Sanger	454 (Roche), Illumina, ABI SOLiD, Ion Torrent	Single-molecule real time (PacBio), Nanopore (Oxford)
Primary distinction	Gold standard	High-throughput via mass parallelisation of sequencing reactions	
			Sequence single molecules, requiring no DNA amplification
Cost	High	Medium	Low
Read length	Long	Short	Very long
Depth	Low	High	Medium
Error rates	Low	Medium	High

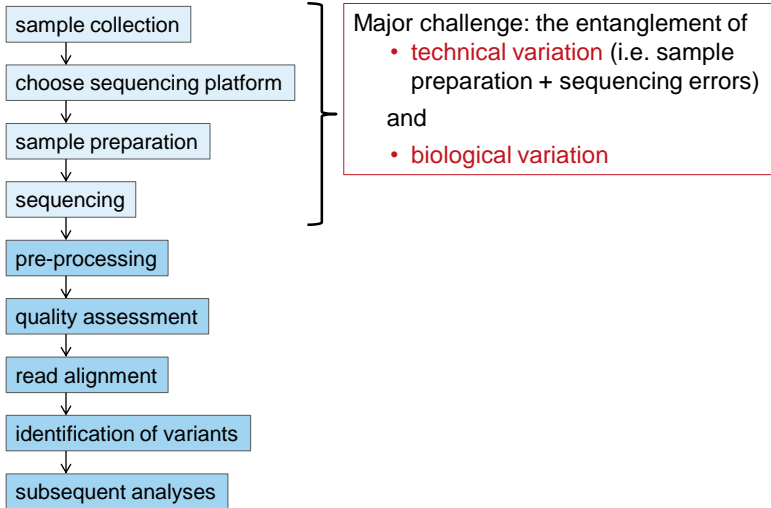
5

Illumina sequencing

- Sequencing-by-synthesis approach.
- Data output (MiSeq reagent kit v3):
 - Read lengths: 2x300bp
 - Number of reads: 44-50 million paired-end reads
- Sequencing error profile:
 - Predominant errors are substitution errors.
 - Average sequencing error rate: ~0.1% per base.
 - Higher error rates (~1% per base) towards sequence ends.
 - Higher error rates on one strand of the pair end reads.

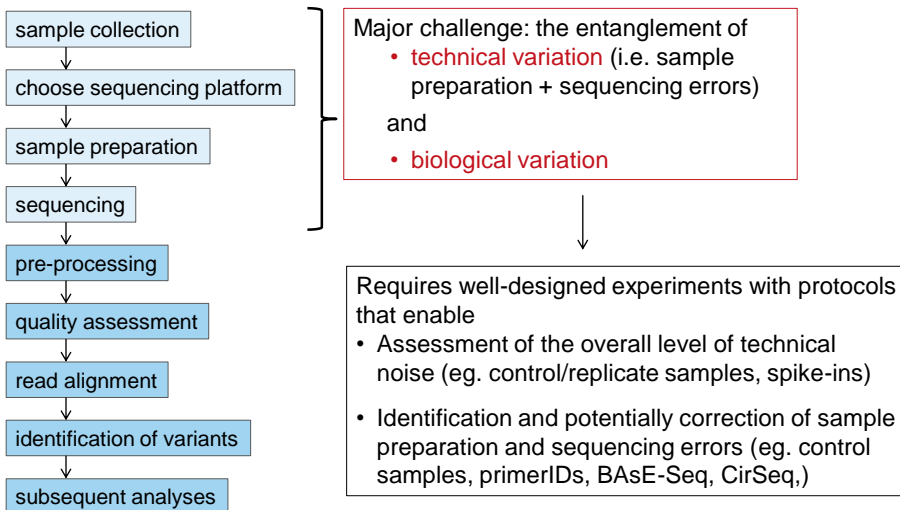
6

Experimental considerations



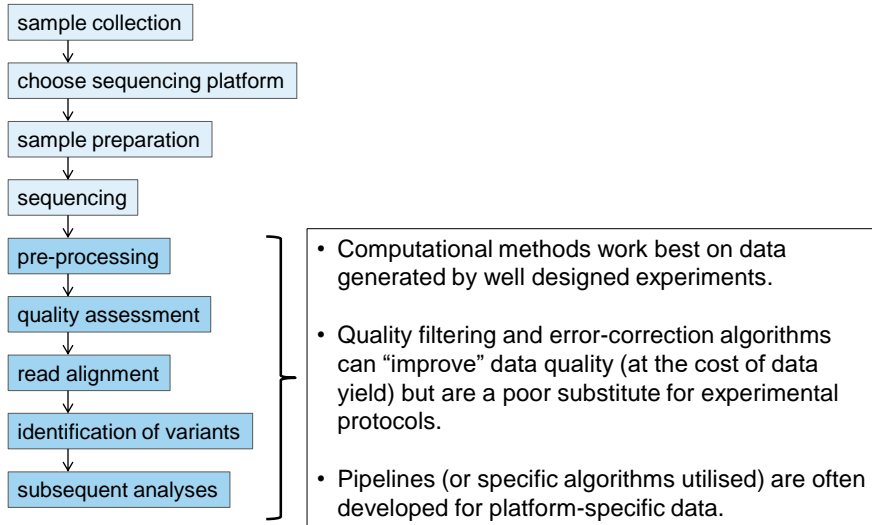
7

Experimental considerations



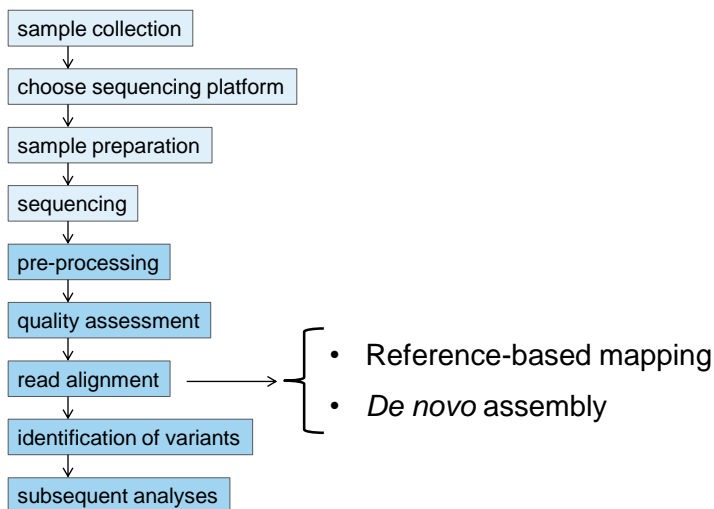
8

Bioinformatics considerations



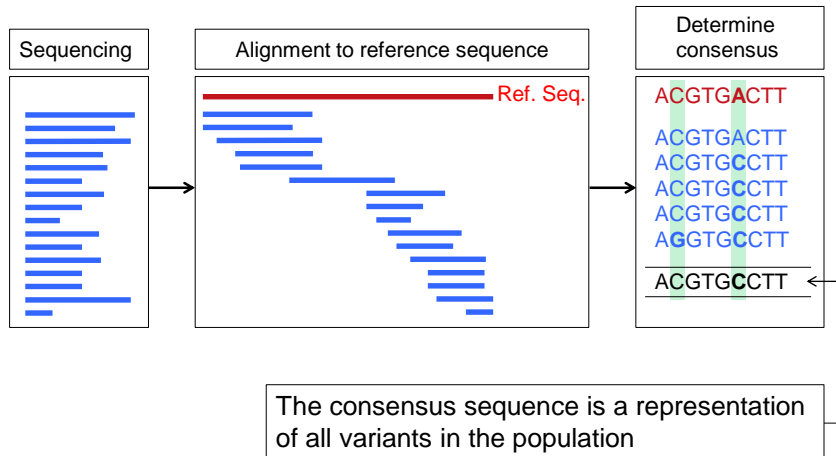
9

Assembly / alignment of reads



10

Reference-based mapping



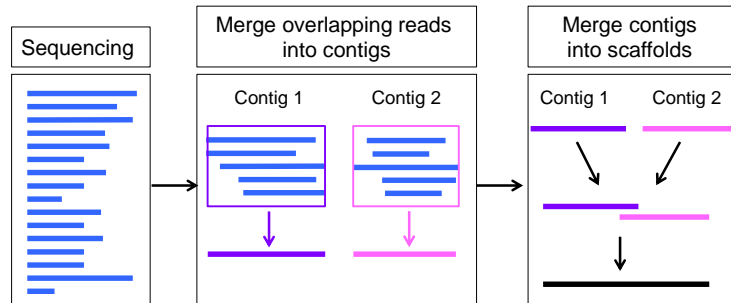
11

Reference-based mapping

- Requires a reference sequence (eg. Los Alamos National Laboratory HIV database, <https://www.hiv.lanl.gov/>).
- Biased towards:
 - Sequences that are more similar to the reference sequence.
 - More abundant variants.
- Reads that substantially differ from the reference sequence align poorly and are often discarded in subsequent analyses.
- Advantage that small variations (eg. SNVs) are more easily positioned.

12

De novo assembly



13

De novo assembly

- No bias towards a reference sequence.
- The assembly is often more fragmented.
- Works better for medium- to large-scale variations.

14

Assembly / alignment of reads

- Assembly using hybrid reference mapping / de novo approaches have some advantages.
- The quality of sequence alignments depends on the strengths and weaknesses of the alignment algorithm used:
 - Computational efficiency.
 - Alignment biases.
Example: Bias introduced in the handling of gaps. Many algorithms do not support gapped alignments, try to minimise gaps at all costs, or preferentially position gaps (eg. in homopolymer regions).
 - Loss of reads.

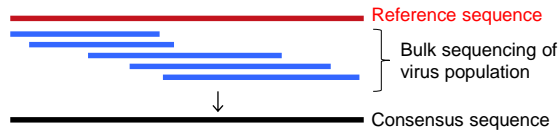
15

Sequencing approaches for HIV diversity

- **Consensus Sanger sequencing**
- **Next-generation sequencing (NGS)**
- **Single genome sequencing (SGS)**

16

Consensus Sanger sequencing

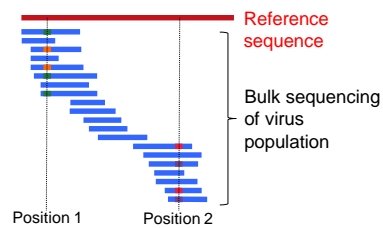


- Has a 20% frequency detection threshold.
- Limited sensitivity to detect low-frequency variants -> underestimation of variant diversity.

17

Next-generation (NGS) sequencing

- High coverage assists in the detection of low-frequency mutations at particular positions, but difficult to link mutations at one position with mutations at another position, due to short read lengths.

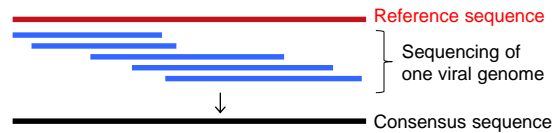


- Coverage can be uneven, with potential fragmentation.
- Higher sequencing error rates -> challenge to distinguish real variation from technical noise (sample preparation + sequencing error).

18

Single genome sequencing

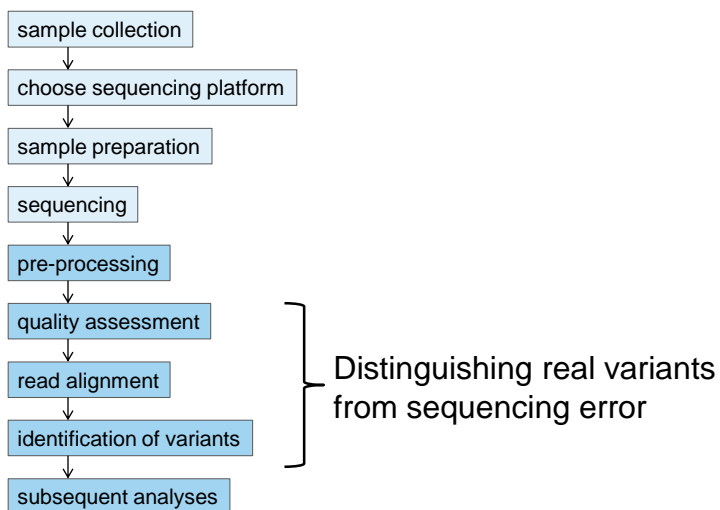
- Isolation of single cell -> nucleic acid extraction -> single genome amplification (SGA) -> sequencing library preparation -> sequencing.



- Advantages:
 - Consensus sequence represents a single viral strain.
 - Enables detection of distant co-occurring mutations.
- Challenges:
 - Efficient physical isolation of individual cells.
 - Minimising amplification bias.

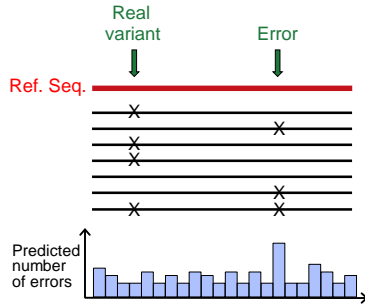
19

Identification of variants



20

Identification of variants



- Identify mismatches with the reference sequence that are more frequently observed than predicted by an error model.
- Various methods for modelling the distribution of errors at each sequence position, based on:
 - Quality scores.
 - Adaptive quality threshold estimated for each sample.
 - Control samples.

21

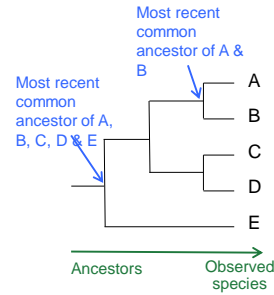
Analysing diversity

- Phylogenetic analysis
- Measures of diversity

22

Phylogenetic analysis

- Provides insights into the viral genetics evolution and relationships between viral variants.
- Assesses the genetic distance between sequences and identifies clusters of similar sequences.
- Various approaches (eg. neighbor joining, maximum likelihood, Bayesian) make different model assumptions and include different parameters.

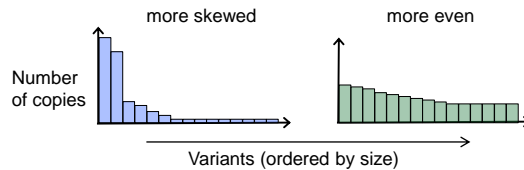


23

Measures of diversity

Counting observed variants

- Highly sensitive to sampling effects, particularly
 - Differences in the copy number distribution



- Differences in sample size (i.e. number of sequences)

24

Measures of diversity

Diversity indices

- Used in ecology and genetics.
- Account for number of unique variants and their copy number distribution, and differences in sample size.
- Compare sample diversity between groups of samples (eg. Simpson's, Shannon, Gini indices)
- Estimate total diversity of population from the samples (eg. Chao estimators)

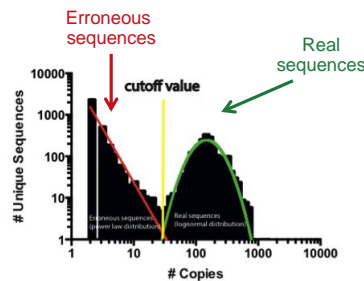
25

Indicators of technical noise contribution

- Ratio of observed variants to input viral templates.



- Over-sequencing copy number distribution.



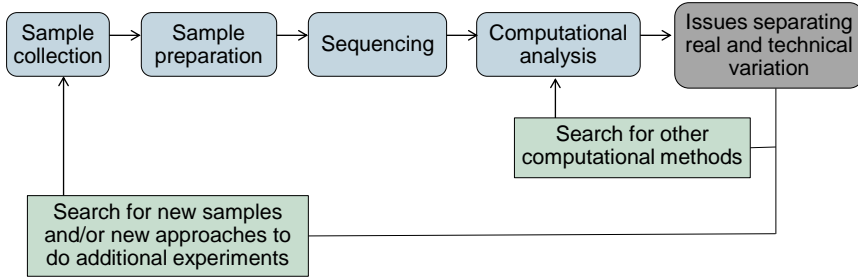
Input: 5000 virions

(Fennessey et al PLoS Path 2017)

26

Concluding remarks

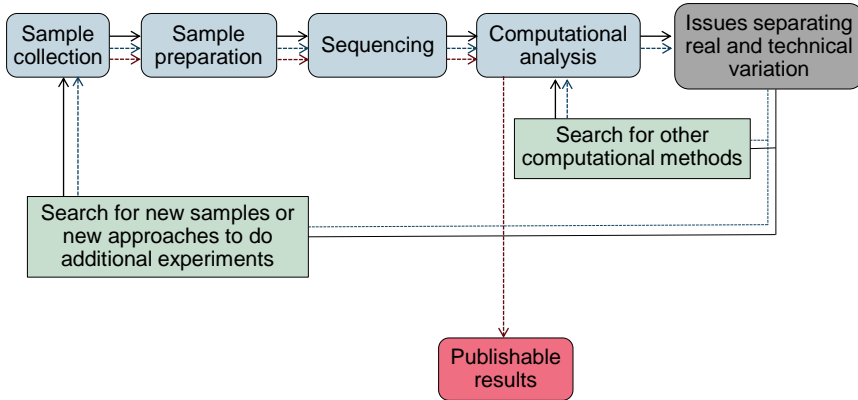
Sequencing project progression often looks like:



27

Concluding remarks

Sequencing project progression often looks like:



28

Concluding remarks

Sequencing project progression should look like:

