# Retroviruses integrate into a shared, non-palindromic motif
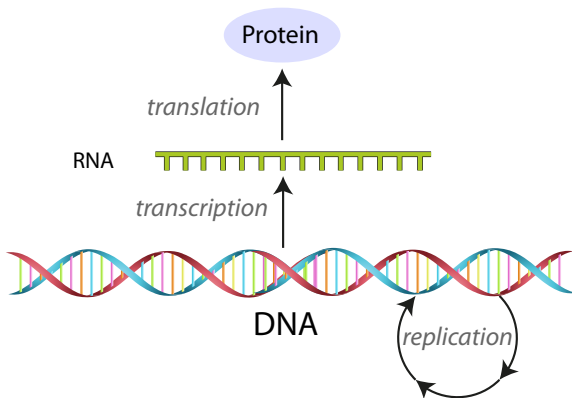
## Paul Kirk

MASAMB 2016, Cambridge

October 4, 2016
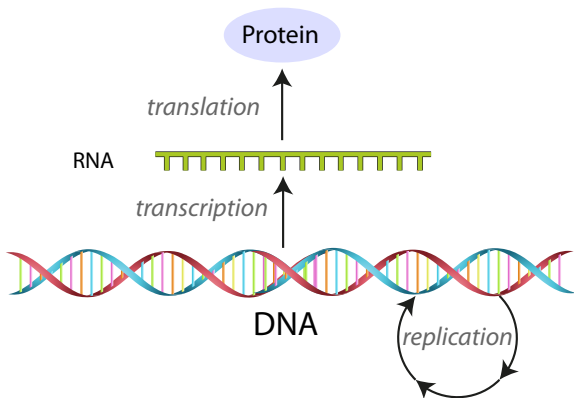
**General** transfers of biological sequential information:

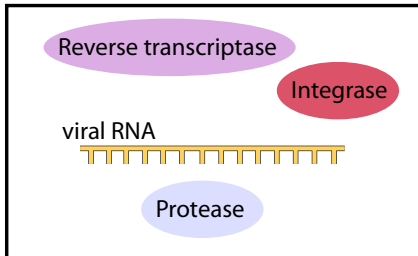**General** transfers of biological sequential information:



There are also **special** transfers of sequential information.

A retrovirus:

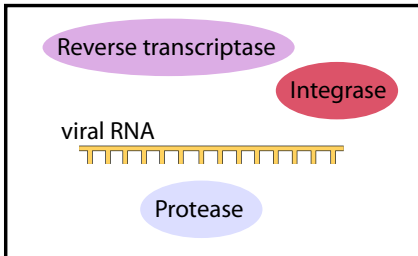Reverse transcriptase

Integrase

viral RNA

Protease
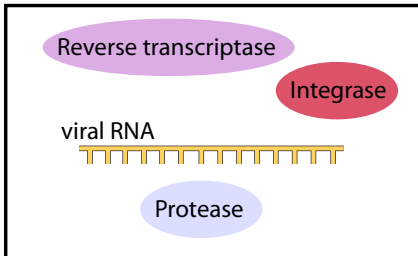
A retrovirus:



Retroviruses are *obligate parasites*: they require a host cell to complete their "life"-cycle.

A retrovirus:



Retroviruses are *obligate parasites*: they require a host cell to complete their "life"-cycle.

Examples: HIV, HTLV-1, . . . .

HOST CELL

host DNA

HOST CELL

# INFECTION

host DNA

HOST CELL

viral RNA

host DNA

HOST CELL

viral RNA

Reverse transcriptase

viral DNA

host DNA

HOST CELL

viral RNA

Reverse transcriptase

Integrase

host DNA  viral DNA  host DNA

provirus

We would like to characterise the target integration site

- i.e. the regions flanking the provirus
- Is there a motif?

Given a collection of integration sites, we can align them according to the position of the provirus. . .

Given a collection of integration sites, we can align them according to the position of the provirus. . .

. . . and then ignore/remove/mask the provirus sequence, so that we just look at the target sites:

# Summarising a collection of target sites

Sequences

**Example
(5 sequences)**

```
...ATC...
...TTA...
...AAC...
...TTC...
...AGC...
```

### Consensus sequence

Just take the most frequent letter at each position: ...ATC...

### Position probability matrix (PPM), *P*

Estimate the probability of each letter at each position:

$$P = \begin{array}{c} A \\ T \\ C \\ G \end{array} \begin{pmatrix} \ldots & 3/5 & 1/5 & 1/5 & \ldots \\ \ldots & 2/5 & 3/5 & 0 & \ldots \\ \ldots & 0 & 0 & 4/5 & \ldots \\ \ldots & 0 & 1/5 & 0 & \ldots \end{pmatrix}$$

## Summarising a collection of target sites

|  | Sequences | Complements | Reverse complements |
|---|---|---|---|
| **Example (5 sequences)** | ...ATC... | ...TAG... | ...GAT... |
| | ...TTA... | ...AAT... | ...TAA... |
| | ...AAC... | ...TTG... | ...GTT... |
| | ...TTC... | ...AAG... | ...GAA... |
| | ...AGC... | ...TCG... | ...GCT... |

### Reverse complement PPM, $P^{(RC)}$

The PPM for the reverse complement sequences:

$$P^{(RC)} = \begin{array}{c} A \\ T \\ C \\ G \end{array} \begin{pmatrix} \ldots & 0 & 3/5 & 2/5 & \ldots \\ \ldots & 1/5 & 1/5 & 3/5 & \ldots \\ \ldots & 0 & 1/5 & 0 & \ldots \\ \ldots & 4/5 & 0 & 0 & \ldots \end{pmatrix}$$

Note: we can get $P^{(RC)}$ from $P$ (and vice versa) by swapping the rows A $\leftrightarrow$ T and C $\leftrightarrow$ G, and reversing the order of the columns.

From 4,521 HTLV-1 target integration sites, we find the consensus:

**AAGTGGATATCCACTT**

From 13,442 HIV-1 target integration sites, we find the consensus:

**TTTGGTAACCAAA**

From 4,521 HTLV-1 target integration sites, we find the consensus:

**AAGTGGATATCCACTT**

From 13,442 HIV-1 target integration sites, we find the consensus:

**TTTGGTAACCAAA**

From 4,521 HTLV-1 target integration sites, we find the consensus:



From 13,442 HIV-1 target integration sites, we find the consensus:

From 4,521 HTLV-1 target integration sites, we find the consensus:



From 13,442 HIV-1 target integration sites, we find the consensus:



The target integration sites are palindromic (as already known!)

# Palindromic PPMs for HTLV-1 and HIV-1 target integration sites

For both HTLV-1 and HIV-1, we have $P^{(RC)} \approx P$

**HTLV-1**



**HIV-1**

# Palindromic sequence logos

**HTLV-1:**



**HIV-1:**

- There is an almost unbelievable amount of symmetry (!)

- There is an almost unbelievable amount of symmetry (!)
- Is this "real"? Do we see evidence of the symmetry within individual sequences, or just at the level of these summaries?

# An attack of aibohphobia

- There is an almost unbelievable amount of symmetry (!)

- Is this "real"? Do we see evidence of the symmetry within individual sequences, or just at the level of these summaries?

- **We introduce a palindrome index to quantify "how palindromic" each sequence is**

# AAGTGGATATCCACTT

$$S = S_{-8}\ S_{-7}\ S_{-6}\ S_{-5}\ S_{-4}\ S_{-3}\ S_{-2}\ S_{-1}\ S_1\ S_2\ S_3\ S_4\ S_5\ S_6\ S_7\ S_8$$

**AAGTGGATATCCACTT**

$$\mathbf{S} = S_{-8}\ S_{-7}\ S_{-6}\ S_{-5}\ S_{-4}\ S_{-3}\ S_{-2}\ S_{-1}\ S_1\ S_2\ S_3\ S_4\ S_5\ S_6\ S_7\ S_8$$

Define

$$\rho(\mathbf{S}) = \frac{1}{n}\sum_{i=1}^{n}\mathbb{I}(s_i = c(s_{-i})),$$

where $2n$ is the sequence length, $\mathbb{I}$ is the indicator function, and $c(x)$ is the complement of $x$ (e.g. $c(T) = A$).

**AAGTGGATATCCACTT**

$$\mathbf{S} = S_{-8}\ S_{-7}\ S_{-6}\ S_{-5}\ S_{-4}\ S_{-3}\ S_{-2}\ S_{-1}\ S_1\ S_2\ S_3\ S_4\ S_5\ S_6\ S_7\ S_8$$

Define

$$\rho(\mathbf{S}) = \frac{1}{n}\sum_{i=1}^{n}\mathbb{I}(s_i = c(s_{-i})),$$

where $2n$ is the sequence length, $\mathbb{I}$ is the indicator function, and $c(x)$ is the complement of $x$ (e.g. $c(T) = A$).

(In practice, we use an "adjusted for chance" version, which is maximally 1, and is 0 if **S** is no more palindromic than expected by chance.)

- The individual sequences are not palindromic

- The individual sequences are not palindromic

- So why do we see palindromes when we average over a large number of sequences?

- One possible explanation is that we have a mix of "forward" and "reverse complement" sequence orientations,

- One possible explanation is that we have a mix of "forward" and "reverse complement" sequence orientations,
  e.g. in the noiseless case

```
Sequence 1:  AATTTAAGTGGAT (Forward)
Sequence 2:  ATCCACTTAAATT (Reverse complement)
Sequence 3:  ATCCACTTAAATT (Reverse complement)
Sequence 4:  AATTTAAGTGGAT (Forward)
Sequence 5:  ATCCACTTAAATT (Forward)
Sequence 6:  AATTTAAGTGGAT (Reverse complement)
```

$$
P = \begin{matrix} A \\ T \\ C \\ G \end{matrix} \begin{pmatrix} 1 & 0.5 & 0 & 0 & 0.5 & 0.5 & 0.5 & 0 & 0.5 & 0.5 & 0.5 & 0.5 & 0 \\ 0 & 0.5 & 0.5 & 0.5 & 0.5 & 0 & 0.5 & 0.5 & 0.5 & 0 & 0 & 0.5 & 1 \\ 0 & 0 & 0.5 & 0.5 & 0 & 0.5 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0.5 & 0 & 0.5 & 0.5 & 0 & 0 \end{pmatrix} = P^{(RC)}
$$

## Analogy

If we have a sample of many real numbers, and we take their mean and find it to be **exactly zero**, one possibility is that this mean is representative of the sample:

## Analogy

If we have a sample of many real numbers, and we take their mean and find it to be **exactly zero**, one possibility is that this mean is representative of the sample:



Another possibility is that we have 2 symmetric components, one positive and one negative:

## Mixture modelling

- We model the sequences as coming from two populations
  - one with PPM $P$; and
  - one with reverse complement PPM $P^{(RC)}$.

$$\pi(S) = \omega\pi(S|P) + (1 - \omega)\pi(S|P^{(RC)}).$$

## Mixture modelling

- We model the sequences as coming from two populations
  - one with PPM $P$; and
  - one with reverse complement PPM $P^{(RC)}$.

$$\pi(S) = \omega\pi(S|P) + (1 - \omega)\pi(S|P^{(RC)}).$$

- Here, $\omega$ is the proportion of sequences coming from the population with PPM $P$.

## Mixture modelling

- We model the sequences as coming from two populations
  - one with PPM $P$; and
  - one with reverse complement PPM $P^{(RC)}$.

$$\pi(S) = \omega\pi(S|P) + (1 - \omega)\pi(S|P^{(RC)}).$$

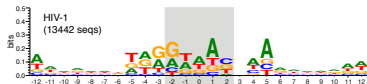- Here, $\omega$ is the proportion of sequences coming from the population with PPM $P$.

- The parameters, $\omega$ and $P$, can be estimated/inferred in numerous ways. I will show results from using an EM-algorithm, but identical results are obtained by: (i) maximum profile likelihood; (ii) Gibbs sampling; (iii) greedy Gibbs.

# Unmixing the forward and reverse sequences



Subpopulation 1 | Subpopulation 2

HTLV-1 (4521 seqs)

HIV-1 (13442 seqs)

- The palindrome is not observed within individual sequences.

- The palindrome is not observed within individual sequences.
- Hypothesis: the palindrome results from a mixture of sequences that contain a non-palindromic motif in approximately equal proportions in "forward" and "reverse complement" orientations

- The palindrome is not observed within individual sequences.

- Hypothesis: the palindrome results from a mixture of sequences that contain a non-palindromic motif in approximately equal proportions in "forward" and "reverse complement" orientations

- Modelling this hypothesis revealed a common nucleotide motif across 4 retroviruses:

    5'-T(N1/2)[C(N0/1)T|(W1/2)C]CW-3'

## Summary

- The palindrome is not observed within individual sequences.

- Hypothesis: the palindrome results from a mixture of sequences that contain a non-palindromic motif in approximately equal proportions in "forward" and "reverse complement" orientations

- Modelling this hypothesis revealed a common nucleotide motif across 4 retroviruses:

    5'-T(N1/2)[C(N0/1)T|(W1/2)C]CW-3'

- Potential implications for understanding retroviral integration.
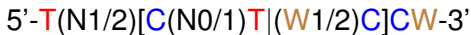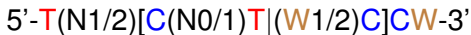
## Summary

- The palindrome is not observed within individual sequences.

- Hypothesis: the palindrome results from a mixture of sequences that contain a non-palindromic motif in approximately equal proportions in "forward" and "reverse complement" orientations

- Modelling this hypothesis revealed a common nucleotide motif across 4 retroviruses:

     5'-T(N1/2)[C(N0/1)T|(W1/2)C]CW-3'

- Potential implications for understanding retroviral integration.

- True validation requires further structural information about retroviral intasomes.

- Accepted for publication in Nature Microbiology.
- Preprint:
  - Kirk, Huvet, Melamed, Maertens & Bangham (2015). Retroviruses integrate into a shared, non-palindromic motif. bioRxiv.

Matlab code (and the HTLV-1 dataset) are available online:

http://www.mrc-bsu.cam.ac.uk/software/
bioinformatics-and-statistical-genomics/

Just click on **retroCode** to download!

## Acknowledgements

Charles Bangham

Maxime Huvet
Anat Melamed
Goedele Maertens

---

Sylvia Richardson
MRC Biostatistics Unit

---

Michael Stumpf
Imperial College Theoretical Systems Biology group

**MRC** | Biostatistics Unit

@pauldwkirk

http://www.mrc-bsu.cam.ac.uk/people/paul-kirk/

# SCIENCE SHOWOFF

**THE CHAOTIC SCIENCE COMEDY CABARET**

with
**STEVE CROSS**
**SARAH BENNETTO**
and loads of Cambridge
science talent

SEASONS
AND A MOVIE

biology
week
2016

UNIVERSITY OF
CAMBRIDGE

Wellcome Trust - Medical Research Council
Cambridge Stem Cell Institute

MONDAY 10th OCTOBER
Portland Arms, Doors 6.30
Tickets £5 from
scienceshowoff.org
or £7 on the door
All ticket money
will go to Parkinson's UK