# Parametric Inference for Recombination in HIV Genomes

Kevin Woods, Oberlin College
(joint with Niko Beerenwinkel and Colin Dewey)

Genomes from two strains of HIV:

<div align="center">

AAAAAA

CCCCCC

</div>

Find new strain:

<div align="center">

AAACCC

</div>

Genomes from two strains of HIV:

AAAAAA

CCCCCC

Find new strain:

AAACCC

Recombination!

**Goal:** Given genomes of a number of parent strains and a new genome, figure out exactly how it is a recombination.

AAA

CCC

New strain:

TAC

**Goal:** Given genomes of a number of parent strains and a new genome, figure out exactly how it is a recombination.

<div align="center">

AAA

CCC

</div>

New strain:

<div align="center">

TAC

</div>

Is it

<div align="center">

TAC ?

</div>

<div align="center">

1 mutation, 1 recombination event

</div>

**Goal:** Given genomes of a number of parent strains and a new genome, figure out exactly how it is a recombination.

<div align="center">

AAA

CCC

</div>

New strain:

<div align="center">

TAC

</div>

Is it

<div align="center">

TAC ?

</div>

<div align="center">

1 mutation, 1 recombination event

</div>

Is it

<div align="center">

TAC ?

</div>

<div align="center">

2 mutations, 0 recombination events

</div>

Tradeoff between mutations and recombinations.

Two parameters $R$ and $M$, the costs of a recombination event or a mutation.

Annotation: a coloring of the new sequence as a recombination of the parent strains.

Given an annotation with $r$ recombination events and $m$ mutations,

$$\text{total cost } = R \cdot r + M \cdot m$$

Minimize total cost over all possible annotations.

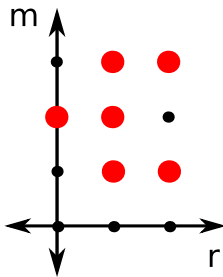Two parameters $R$ and $M$, the costs of a recombination event or a mutation.

Annotation: a coloring of the new sequence as a recombination of the parent strains.

Given an annotation with $r$ recombination events and $m$ mutations,

$$\text{total cost } = R \cdot r + M \cdot m$$

Minimize total cost over all possible annotations.

But what are $R$ and $M$?

Parents:

AAA

CCC

New strain:

TAC

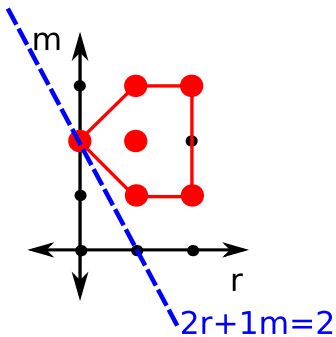8 possible annotations.

Graph $(r, m)$ for each
annotation.

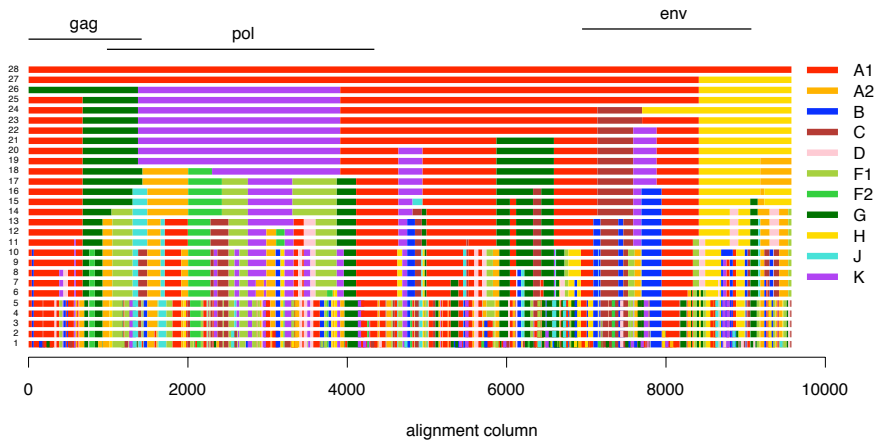Example: $R = 2$, $M = 1$

Maximize

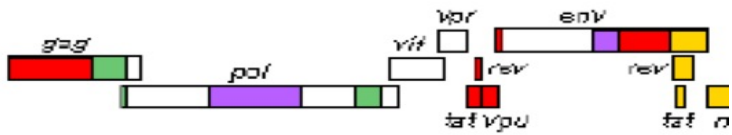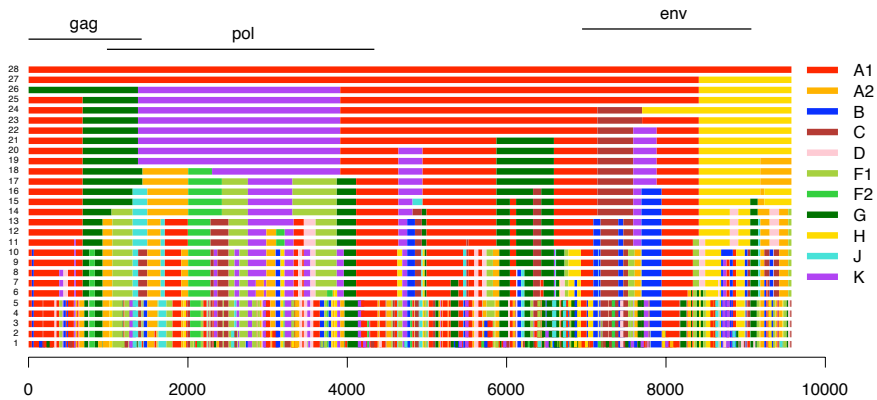$$2r + 1m$$

over all points.

Linear Programming!



m

r

2r+1m=2

The vertices of this polytope are all of the information we need to understand this problem for all parameters (Pachter, Sturmfels).

We can use all of our tools and algorithms from polytopal combinatorics.

This approach gives us a manageable number of interesting annotations.

alignment column

This can be translated to a probabilistic Hidden Markov Model, detailing

$P(\text{annotation}, \text{child sequence} \mid \text{parameters}, \text{parent sequences})$

For fixed parameters "minimizing cost" is finding annotation that maximizes this probability (MAP estimate).

What are the parameters?

What are the parameters?

Classical answer: Maximum Likelihood Estimate: Find the parameters that maximize

$$P(\text{child sequence} \mid \text{parameters}, \text{parent sequences})$$

What are the parameters?

Classical answer: Maximum Likelihood Estimate: Find the parameters that maximize

$$P(\text{child sequence} \mid \text{parameters, parent sequences})$$

Our approach allows interesting alternative: Given prior distribution on $P(\text{parameters})$ (uniform?), compute, for each vertex,

$$P(\text{vertex is MAP estimate} \mid \text{parent sequences, child sequence})$$

This approach gives us a manageable number of interesting annotations, each with a numerical score associated with them.

alignment column

og log posterior