

Pair HMM probabilistic realignment in HaplotypeCaller and Mutect

David Benjamin^{*†}

Broad Institute, 75 Ames Street, Cambridge, MA 02142

(Dated: August 10, 2017)

After generating candidate haplotypes, the GATK tools **HaplotypeCaller** and **Mutect** realign reads against these haplotypes to obtain a matrix of likelihoods for each read to be derived from each haplotype. Here we describe the probabilistic model specifying this likelihood as well as its computational implementation. We do not describe the translation of this implementation into native code optimized for vectorized architectures.

I. THE PAIR HMM MODEL

We want to calculate the probability $P(\mathcal{R}|\mathcal{H})$ of read \mathcal{R} to be sequenced from haplotype \mathcal{H} , where the haplotypes are sufficiently long that reads are contained within them. This likelihood is the sum of likelihoods of all possible alignments \mathcal{A} of \mathcal{R} to \mathcal{H} :

$$P(\mathcal{R}|\mathcal{H}) = \sum_{\mathcal{A}} P(\mathcal{R}, \mathcal{A}|\mathcal{H}) \quad (1)$$

We represent alignments as sequences

$$\text{alignment} = \{(i_1, j_1, s_1), (i_2, j_2, s_2) \dots (i_N, j_N, s_N)\}, \quad (2)$$

where i and j represent positions within the read and haplotype, respectively, and $s_n \in \{M, I, D\}$ represents the states of match, insertion, and deletion of the read relative to the haplotype. For example, an alignment $\{(1, 10, M), (2, 11, M), (3, 12, M), (4, 13, M), (4, 14, D), (4, 15, D), (5, 16, M), (6, 17, M), (7, 17, I), (8, 18, M)\}$ means that positions 1 - 4 of the read match positions 10 - 13 of the haplotype, followed by a two-base deletion (advancing from 13 to 15 in the haplotype without advancing in the read), followed by a match of read positions 5 - 6 with haplotype positions 16 - 17, followed by an insertion at read position 7, followed by a match. The allowable transitions $(i_n, j_n) \rightarrow (i_{n+1}, j_{n+1})$ are

- $(i, j, M/D/I) \rightarrow (i+1, j+1, M)$: match of read position $i+1$ with haplotype position $j+1$
- $(i, j, M/D) \rightarrow (i, j+1, D)$: deletion after read position i – haplotype position $j+1$ is deleted
- $(i, j, M/I) \rightarrow (i+1, j, I)$: insertion at read position $i+1$

Note that the state label s seems redundant because it can be reconstructed from the sequence of i and j . While this is true, our model treats indel starts differently from indel continuations, and by including the state label we can distinguish these conveniently¹.

The read-alignment likelihood has two components. First is the probability of the sequence of match, insertion, and deletion states, which is

$$P(\mathcal{A}) = \prod_k T_{s_k, s_{k+1}}, \quad (3)$$

where T is a constant matrix of state transition probabilities². Next is the emission probability of the read bases

^{*} The author took no part in development of the methods described below – credit belongs to several others on the GATK team.

[†]Electronic address: davidben@broadinstitute.org

¹ That is, we can tell which type of transition it is by looking back one unit instead of two. Thus we have a first-order Markov model instead of a second-order Markov model.

² The gap continuation probability corresponds to a phred-scaled quality that is set by the **gcPHMM** argument, which is 10 by default, implying that $T_{DD} = T_{II} = 0.1$ and $T_{DM} = T_{IM} = 0.9$. The other transitions are hard-coded based on a phred-scaled indel start quality of 45. That is, $T_{MD} = T_{MI} = 10^{-4.5}$ and $T_{MM} = 1 - 2 \times 10^{-4.5}$. This fact is somewhat obscured in the code, which gives a position index to the transition matrices. Nonetheless, the array representing transitions is constant over this index.

given the haplotype bases they align to and the base qualities:

$$P(\mathcal{R}|\mathcal{A}, \mathcal{H}) = \prod_k P(r_{i_k}|h_{j_k}, q_{i_k})^{\mathbb{I}[s_k=M]}, \quad (4)$$

where r_m and h_n are the m th read bases and n th haplotype base and q_m is the quality of the read's m th base. Note that this only includes alignments in the match state. The per-base emission is given directly from the definition of base quality:

$$P(b_2|b_1, q) = \begin{cases} \epsilon(q)/3 & (b_1 \neq b_2) \\ 1 - \epsilon(q) & (b_1 = b_2) \end{cases} \quad (5)$$

II. DYNAMIC PROGRAMMING

Define the matrices M , I and D by M_{ij} = the total likelihood of *all* paths from the beginning of the read to position i that end in a match state, and likewise for D and I . Then the recursions

$$M_{ij} = P(r_i|h_j, q_i) (M_{i-1,j-1}T_{MM} + I_{i-1,j-1}T_{IM} + D_{i-1,j-1}T_{DM}) \quad (6)$$

$$I_{ij} = M_{i-1,j}T_{MI} + I_{i-1,j}T_{II} \quad (7)$$

$$D_{ij} = M_{i,j-1}T_{MD} + D_{i,j-1}T_{DD} \quad (8)$$

define the entire pair HMM algorithm:

Algorithm 1 Pair HMM algorithm

```

1: Initialize  $M_{0,j} = I_{0,j} = 0$  and  $D_{0,j} = 2^{1020}/|\mathcal{H}|$  for  $1 \leq j \leq |\mathcal{H}|$ .
2: for  $1 \leq i \leq |\mathcal{R}|$  do
3:   for  $1 \leq j \leq |\mathcal{H}|$  do
4:     Calculate  $M_{ij}$  via Equation 6.
5:     Calculate  $I_{ij}$  via Equation 7.
6:     Calculate  $D_{ij}$  via Equation 8.
7:   end for
8: end for
9: Total likelihood  $P(\mathcal{R}|\mathcal{H})$  is  $\sum_j (M_{\mathcal{R},j} + I_{\mathcal{R},j})$ .
```

That the $i = 0$ rows of M and I are initialized to zero corresponds to starting at an imaginary position one base before the read start in a deletion state. The factor of $1/|\mathcal{H}|$ corresponds to a flat prior on which j an alignment starts at. The curious factor of 2^{1020} is a huge number to prevent underflow – all multiplications are by numbers less than 1, so we needn't worry about overflow – which is a much more efficient approach than performing the computation in log space. The omission of D from the returned value recognizes the fact that a terminal deletion is meaningless.

Finally, we note a shortcut that the GATK exploits: when two consecutive haplotypes agree up to the k th position, the first k columns of M , D and I are recycled and the inner loop is over $k < j \leq \mathcal{H}$.