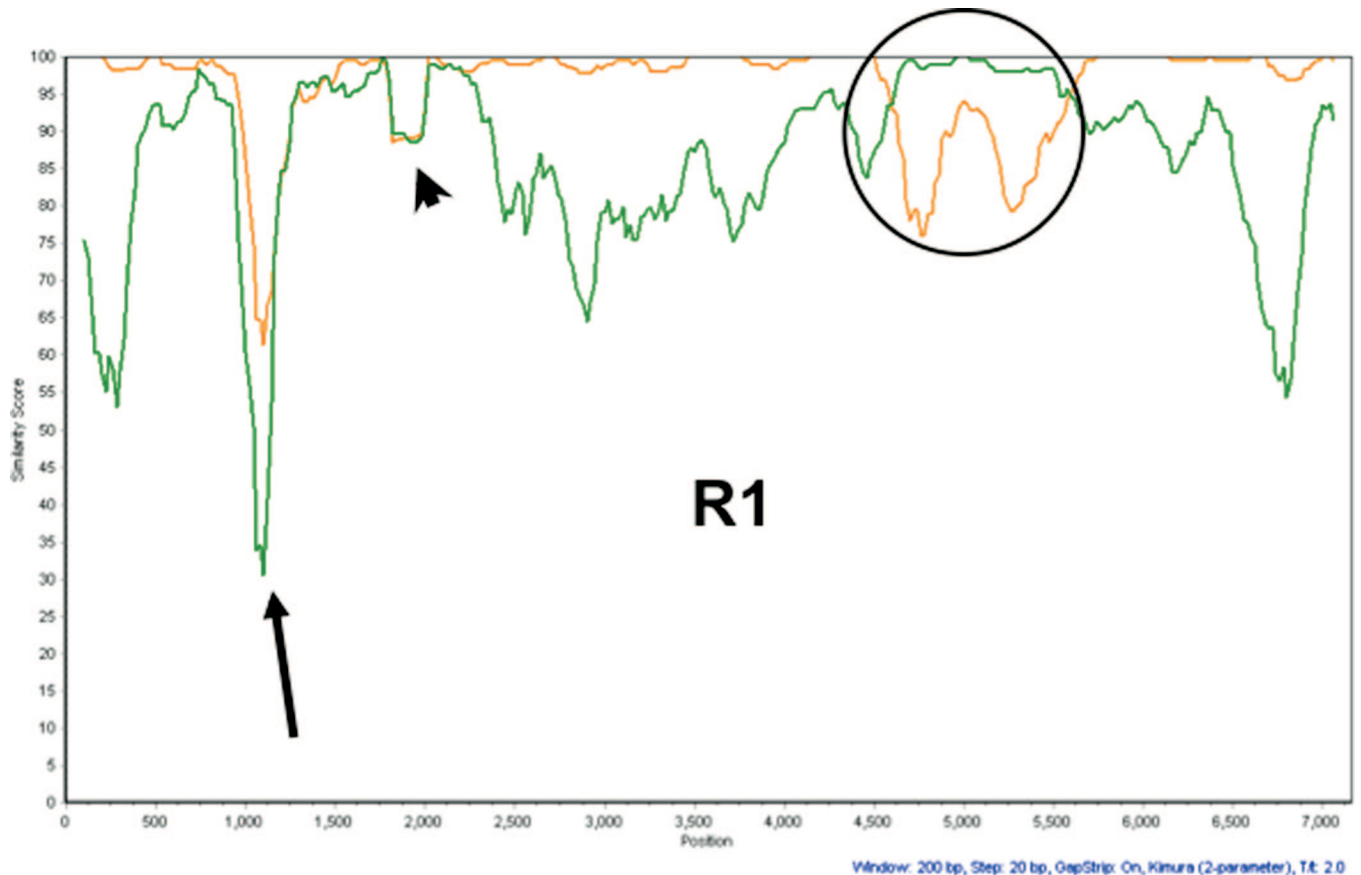


# Supporting Information

Sharma *et al.* 10.1073/pnas.0805694105



**Fig. S1.** SimPlots illustrate the similarity of recombinants R1 to R5 to either parental element CRM1A or CRM1B in different regions of the element. BioEdit was used to generate an 80% consensus sequence for each CRM1 subgroup using separate multiple sequence alignments generated by ClustalX from each of the two parental (CRM1A and CRM1B) and five recombinant (R1–R5) subgroups. The similarity of each recombinant to CRM1A (green) and CRM1B (orange) was plotted using SimPlot with a sliding window of 200 nt and a step size of 20 nt. Detailed recombination breakpoints for all recombinants are shown in Fig. S2. (a) Recombinant R1 more closely resembles parent CRM1B along most of its length, with the exception of the RNaseH region (circled). The R1 consensus was derived from all 89 R1 sequences. Note that the sequence similarity is close to 100% to one or the other parent, except for the region around 1,100 nt (arrow) and 1,900 nt (arrowhead). The former is a highly variable AT-rich region of the UTR containing a duplication within the R1 consensus that results in a big gap in the alignment of A and B with R1, while the latter is a UTR region with high sequence variation. High sequence variation among individual B sequences in part of the UTR indicates that the consensus B sequence is not representative of the B parent in this region. (b) The R1 consensus for this figure was derived from three of the four oldest R1 sequences. One R1 sequence was not included because of the presence of a duplication in the UTR. Note that the consensus derived from these oldest R1 elements shows high sequence similarity with “B” because they lack the duplication (arrow) and the ambiguous region of the UTR (arrowhead), resulting in a much cleaner (and slightly shorter) SimPlot (compare with Fig. S1 a). (c–f) SimPlots of recombinants R2 to R5 demonstrate sequence similarities of each recombinant to either parent CRM1A or CRM1B in different regions of the element, as represented schematically in Fig. 1.

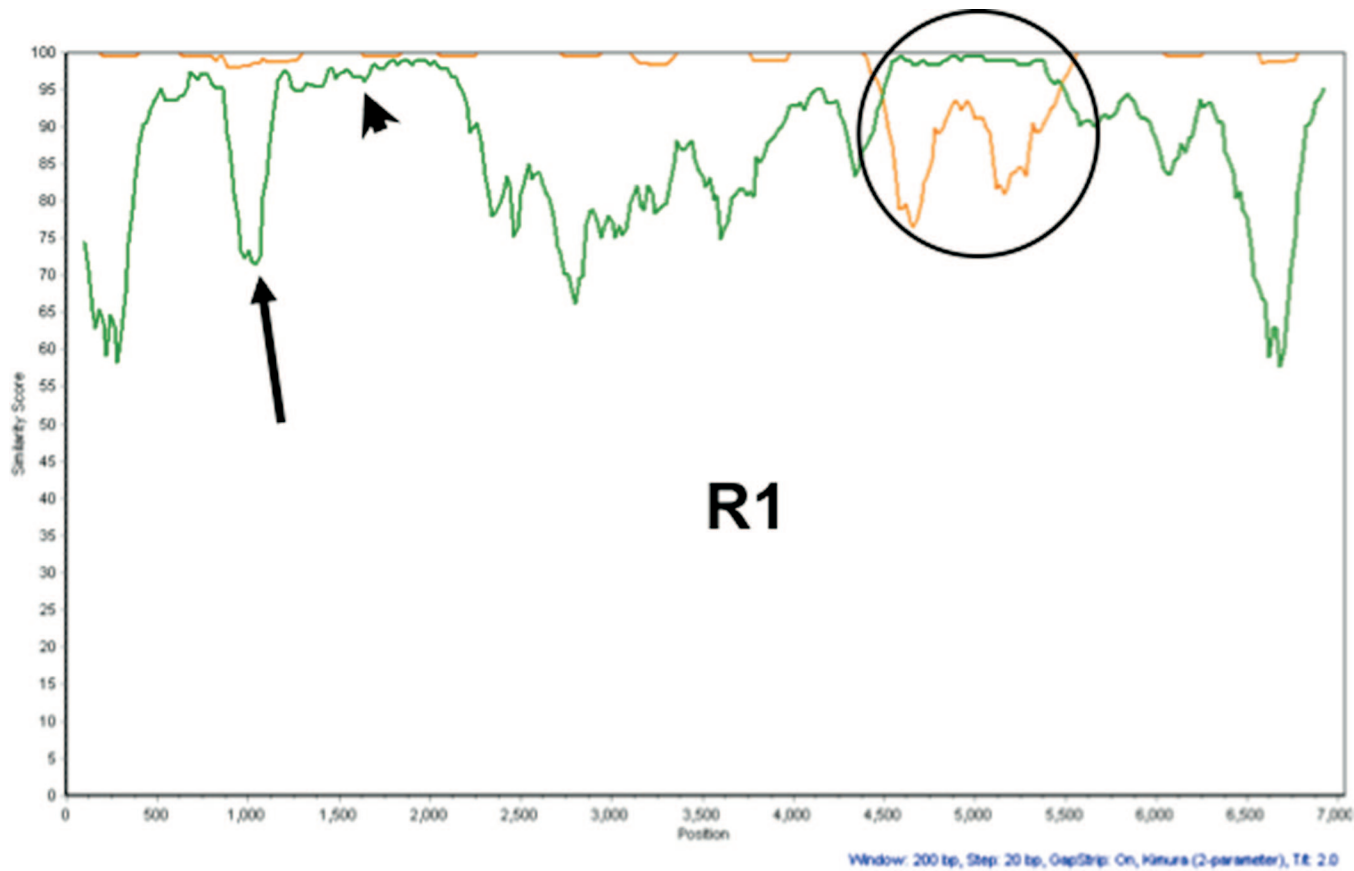


Figure S1. Continued



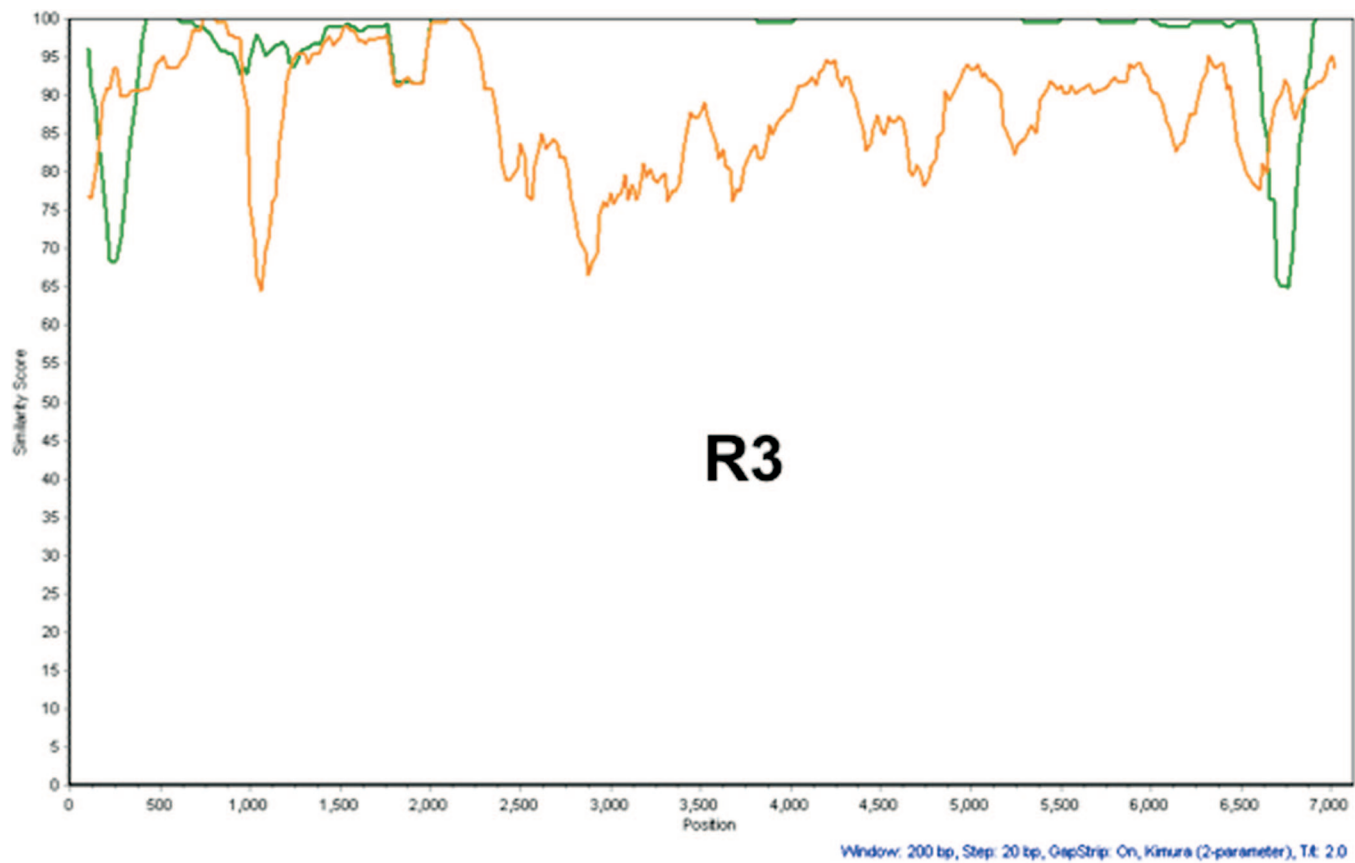


Figure S1. Continued



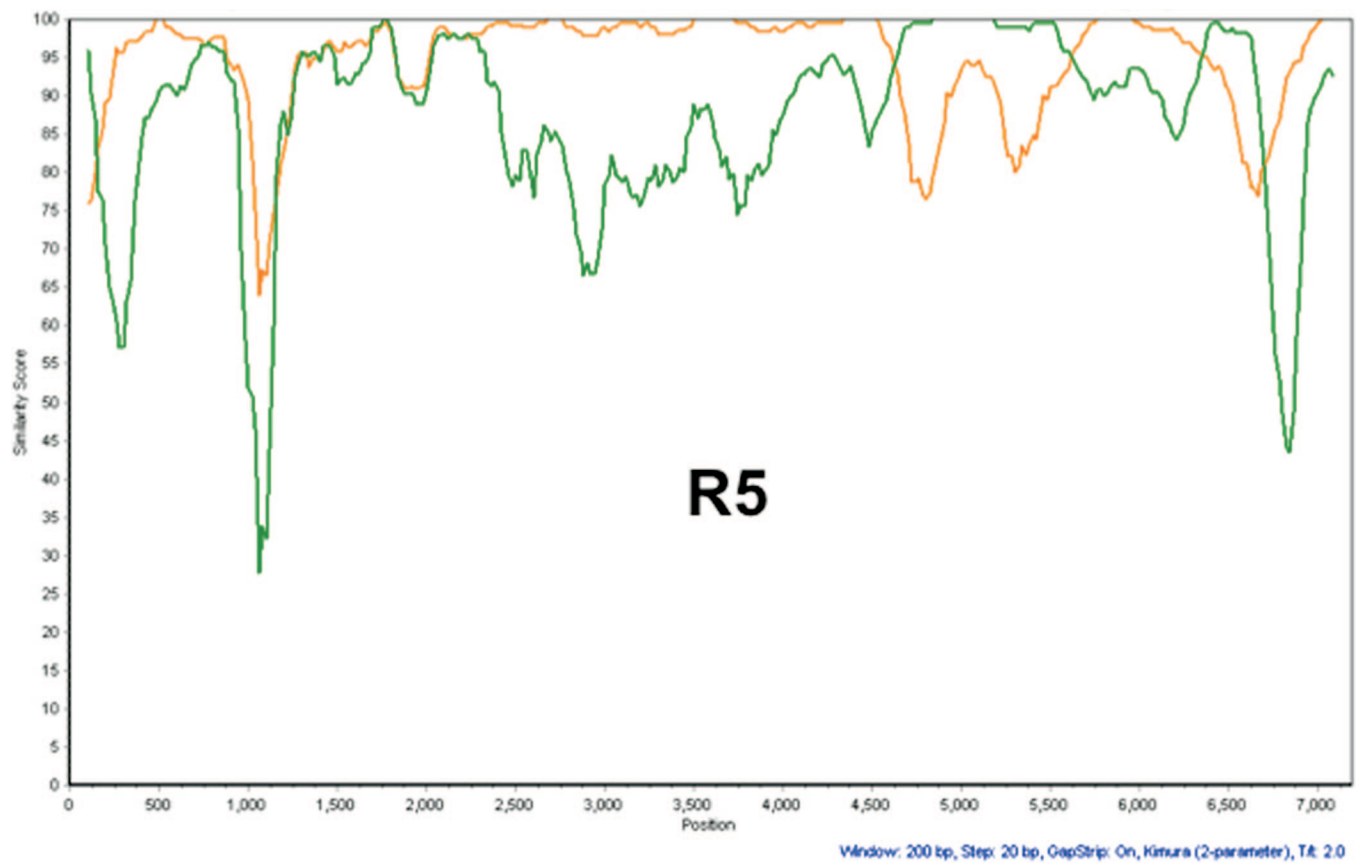
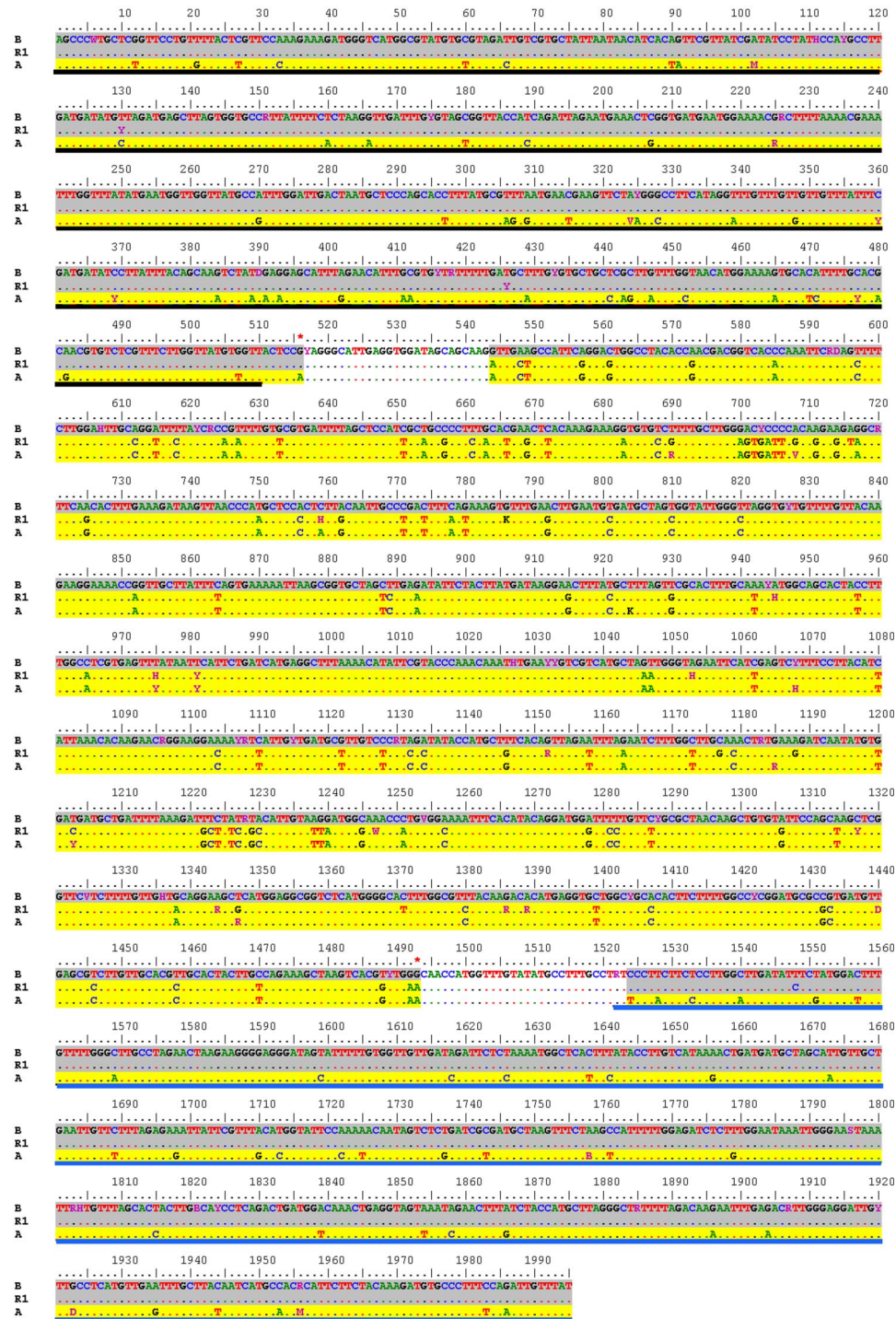


Figure S1. Continued



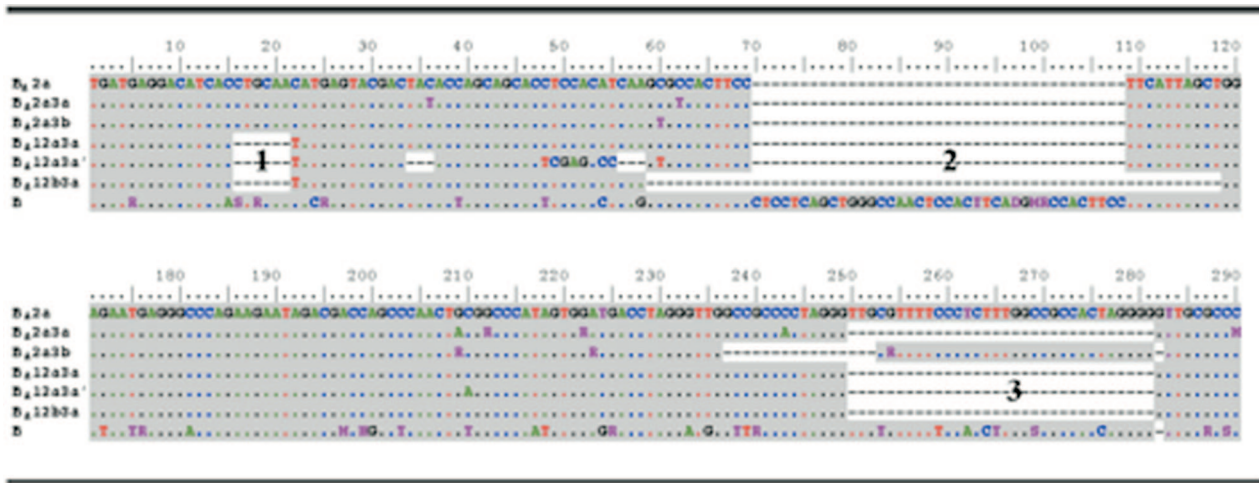
Supporting Figure 2a:



**Fig. S2.** Recombination breakpoints in CRM1 and CRM4 elements. Multiple sequence alignments of the recombinant CRM1 and CRM4 subgroups with their respective parents are shown. The consensus sequence of the recombinant is shown in the middle and those of the parents are at the top and bottom of each alignment. Identities to the top-most sequence are indicated by (.) dots. The recombinant sequence derived from A and B lineages is highlighted in yellow and gray, respectively. The recombination breakpoints identified by “2 parental, 1 derived” Maximum Chi Squared Test is marked by asterisk “\*” for each recombinant. Note: The A polyprotein consensus was generated using the A polyproteins associated with the A and R3 LTRs, while that for the B polyprotein was generated from B polyprotein associated with the B LTRs and their deletion derivatives. The  $B_{\Delta}$  indicates subgroup  $B_{\Delta}12a3a$ , while R1, R2, R3, R4 and R5 represent different recombinant groups. (a) Two recombination breakpoints flank the  $RH_A$  in CRM1 R1 (middle sequence). The flanking RT and Integrase domains (black and blue underline, respectively) are included in this alignment to better visualize the recombination breakpoints. Both recombination breakpoints map upstream of a region of high sequence similarity shared by CRM1 A and B. (b–e) The recombination breakpoints determined for CRM1 recombinants R2, R3 and R5, as well as CRM4 R1.







**Fig. S3.** Deletion derivatives of B LTR. Each deletion derivative of the CRM1B LTR is defined by a characteristic deletion as indicated on this multiple sequence alignment of the 5' end of the CRM1B LTR and its deletion derivatives. The deletions are numbered (1, 2, or 3) based on the region with the deletion, followed by a letter to indicate different extent of deletion. B<sub>Δ12a3a'</sub> differs from B<sub>Δ12a3a</sub> by two additional trinucleotide deletions and base changes.





## Putative Recombination Breakpoints in the Opie Polyprotein region

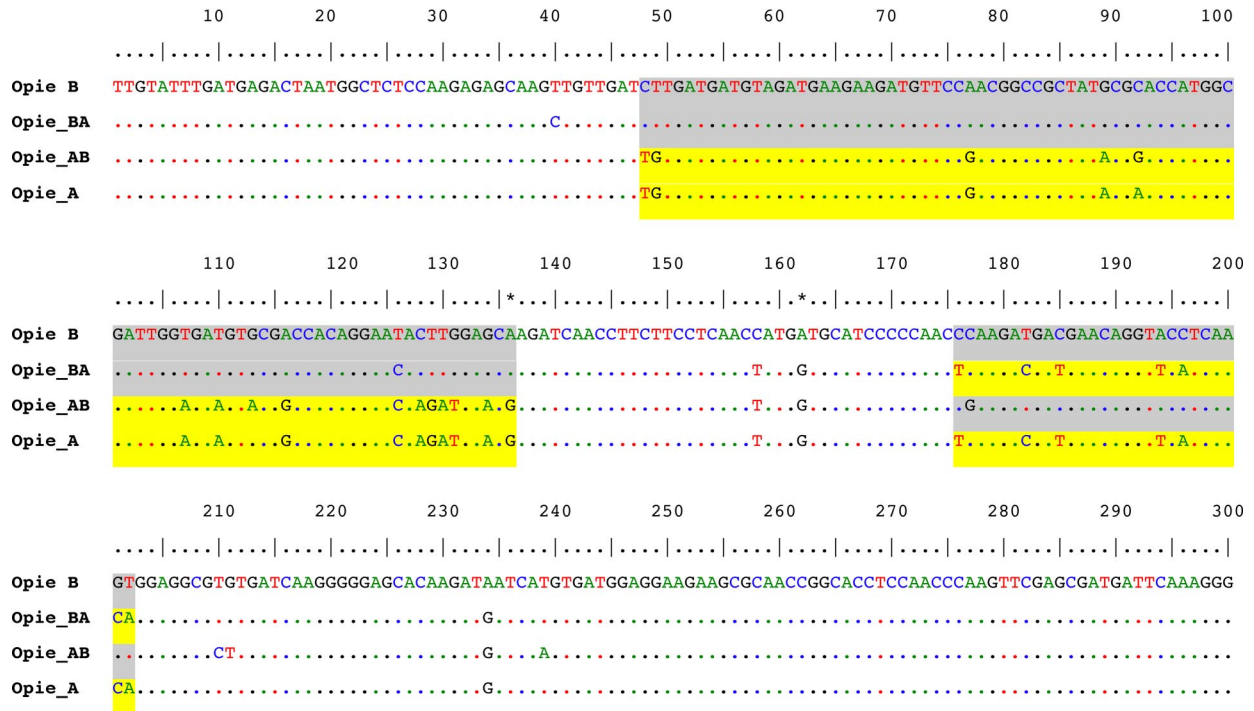


Fig. S6. Recombination breakpoints in Opie polyprotein. Multiple sequence alignments of two recombinant Opie subgroups with two “parents” are shown. The consensus sequence of the recombinants is shown in the middle and those of the parents are at the top and bottom of each alignment. Note that in these highly recombined elements, the parental and recombinant designations are arbitrary. Identities to the top-most sequence are indicated by (.) dots. The sequence derived from A and B lineages is highlighted yellow and gray, respectively. The recombination breakpoints identified by “2 parental, 1 derived” Maximum Chi Squared Test is marked by asterisk for each recombinant (after 136 for AB and after 162 for BA).