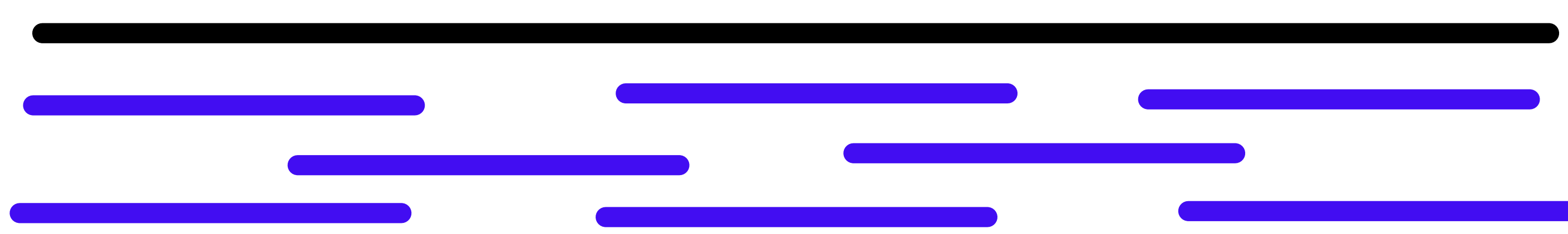


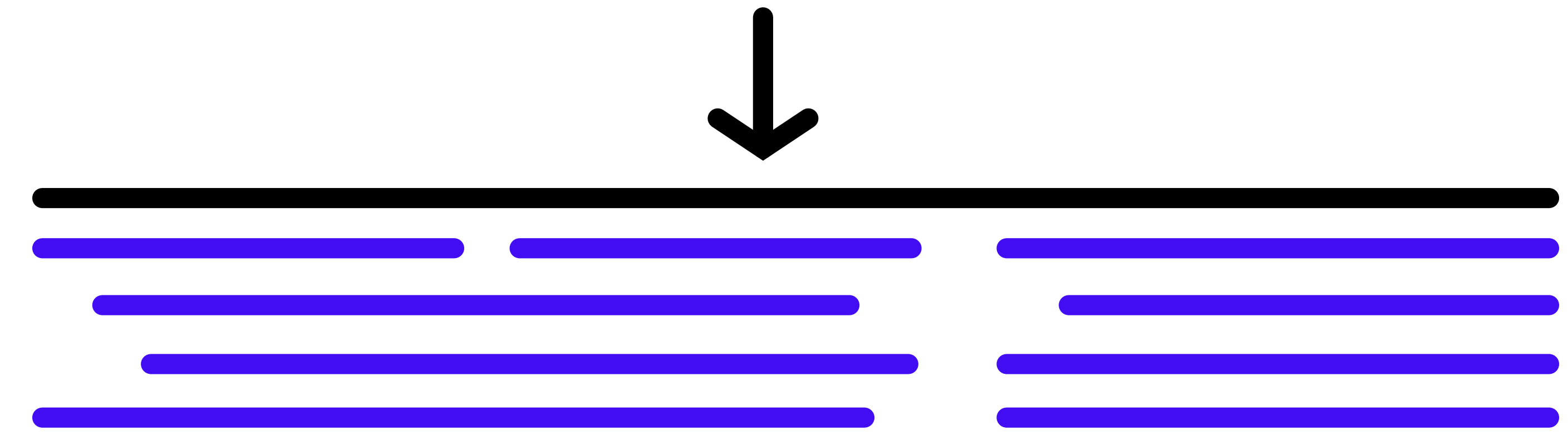
# Protein-Coding Sequence Assembly

## Assembly - Round 1



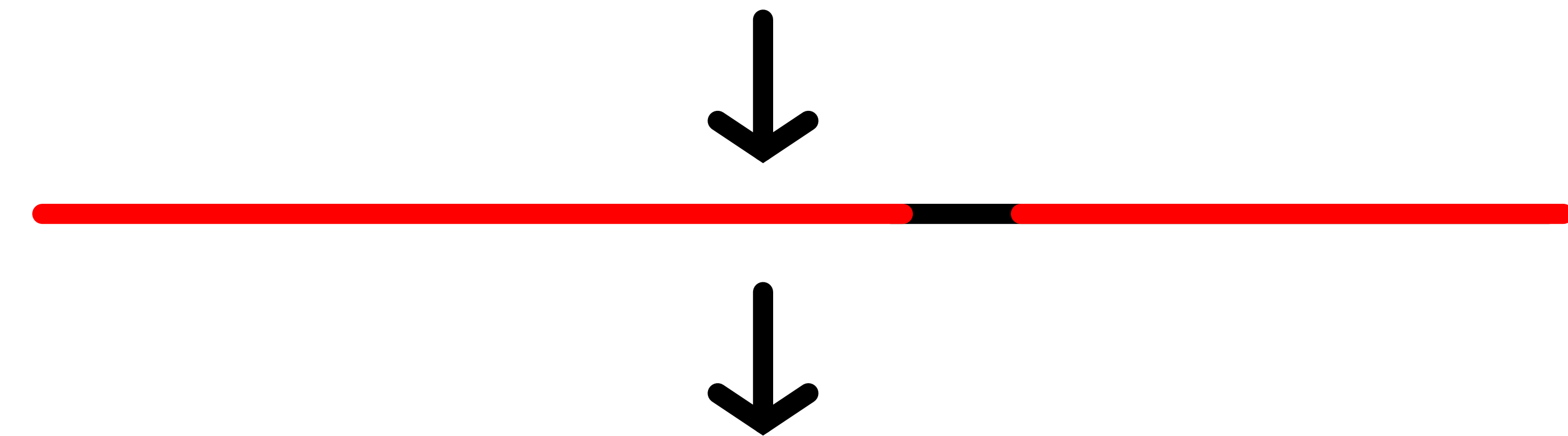
### *Oryzias latipes* reference CDS

Align raw reads to reference sequence using BWA-MEM.



### Multiple sequence alignment (gapped)

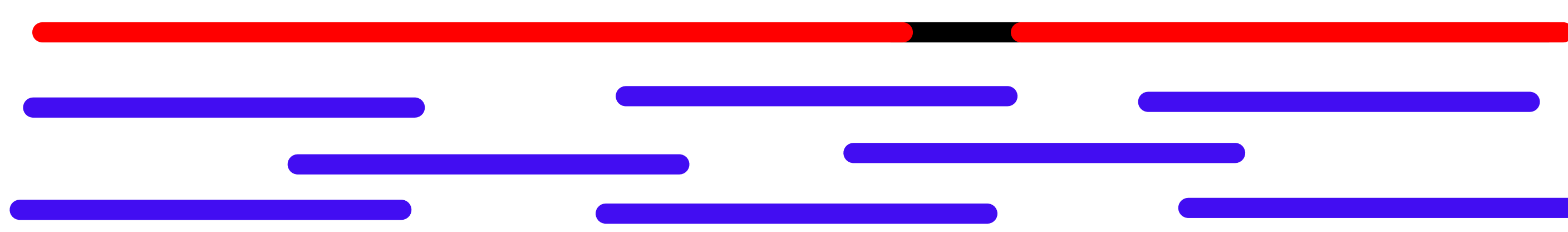
Sp. 1 Generate consensus sequences for each species.  
Sp. 2 Gaps present due to poor assembly in some regions of the CDS.  
Sp. 3  
Sp. 4



### Generate new reference CDS

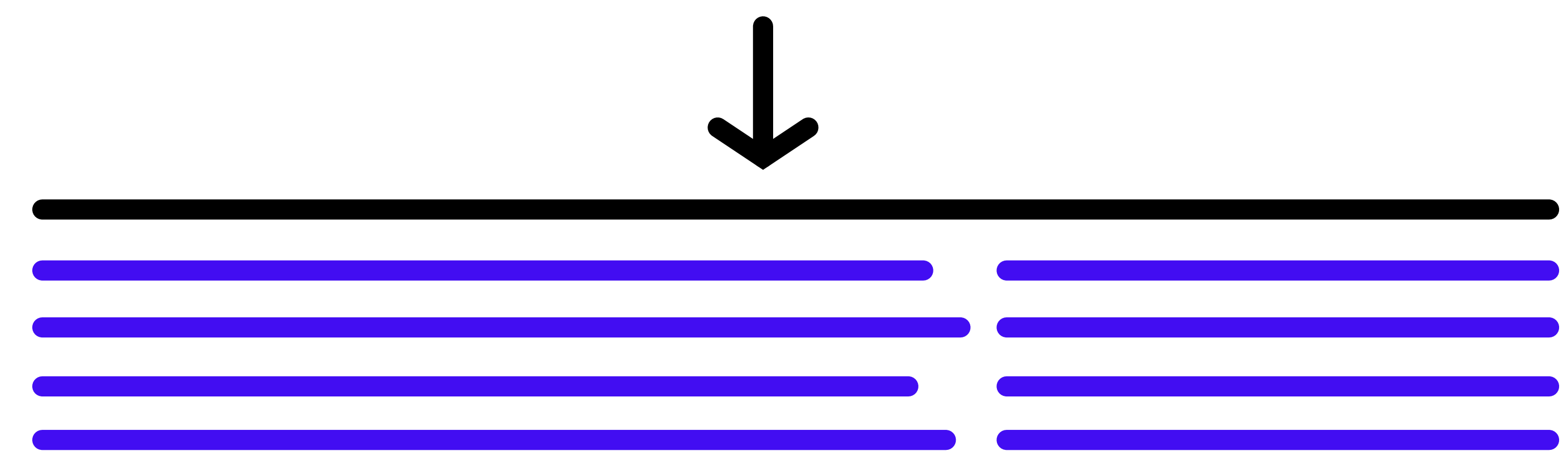
New reference sequence is based on most complete CDS assembled for a species. Gaps are filled in with sequence from *Oryzias latipes* to avoid introducing gaps into the new alignment.

## Assembly - Round 2



### Align raw reads to new CDS

Align raw reads to new "makeshift" reference sequence.



### Multiple sequence alignment (less gapped)

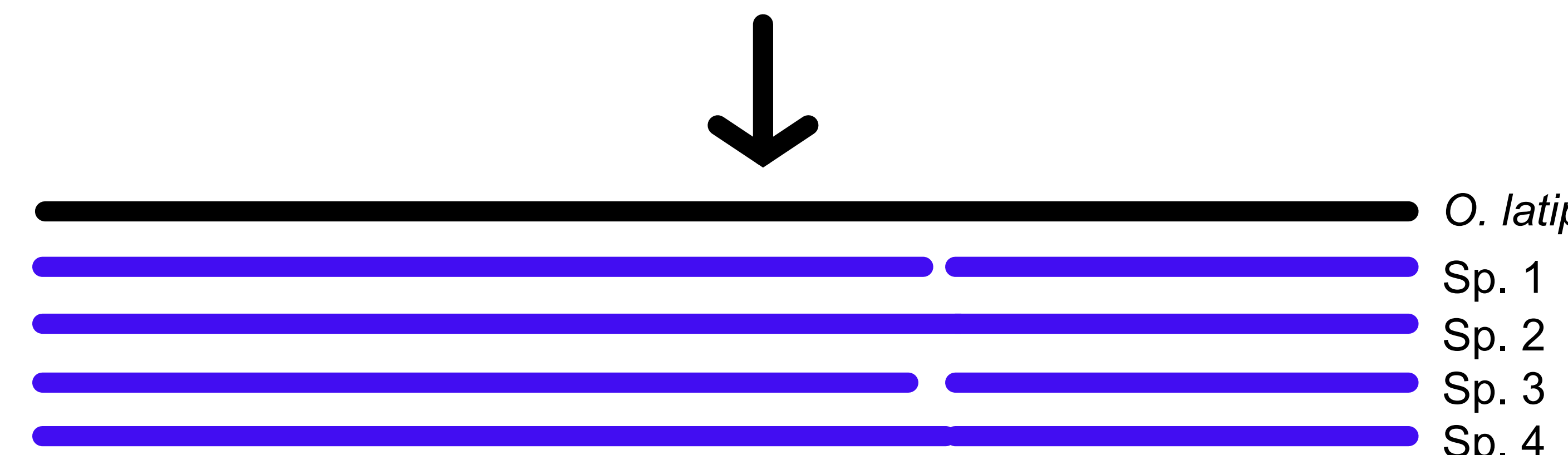
*O. latipes*  
Sp. 1  
Sp. 2 New MSA is less gapped, the new reference sequence is replaced by the original CDS of *Oryzias latipes* to obtain a complete MSA  
Sp. 3  
Sp. 4

## Alignment Cleaning for PAML



### Clean multiple sequence alignment

*O. latipes*  
Sp. 1 If 30%+ of a codon position has gaps, or 10%+ of  
Sp. 2 premature stop codons, those sites are removed  
Sp. 3 to ensure gaps/stop codons do not interfere with  
Sp. 4 PAML analyses.



### Final multiple sequence alignment

*O. latipes*  
Sp. 1 Final, edited MSA has fewer gaps for each species.  
Sp. 2  
Sp. 3  
Sp. 4