

The Problem: Selecting From Competing Annotations

- A common annotation process:
- Compare an un-annotated sequence (the *Target*) to a large collection of known sequences (the *Queries*)
  - Sometimes, more than one of the Queries shows a significant match to the Target
  - Pick ‘true’ matching query based on highest alignment score
  - Falsely implies certainty of the true match
  - When two or more sequences match with high alignment scores, this method is no longer reliable, because it is possible that either one is the true sequence.
  - Relevant in databases with highly similar sequences



Reliability of Competing Annotations

Annotations of biological replicates should agree. Among transposable elements in humans, more than 10% of annotations disagree.

Confidence from Alignment Scores

$Q : q_1, q_2, q_3, \dots, q_n$  (where  $q$  can be a sequence or pHMM)

Define  $P(q_i|t)$  : probability that the true label of  $t$  is  $q_i$

Assuming alignments cover all possible explanations for  $t$ :

$$Conf(q_i|t) = \frac{P(q_i|t)}{\sum_j P(q_j|t)}$$

We don't have  $P(q_i|t)$ , but Bayes' rule says:

$$P(q_i|t) = P(t|q_i) * \frac{P(q_i)}{P(t)}$$

Assuming uniform distribution over  $P(q_i)$ , and with a fixed  $t$ ,  $P(q_i|t) \propto P(t|q_i)$ . So,

$$Conf(q_i|t) = \frac{P(t|q_i)}{\sum_j P(t|q_j)}$$

Probabilities ↔ Scores

Given scoring matrix, alignment scores are (scaled) log odds ratios:

$$Score(a, b) = \text{int} \left( \lambda * \log_2 \frac{P(a, b | \text{homology})}{P(a)P(b)} \right)$$

$\lambda$  : scaling factor of scoring matrix

Over entire alignment:

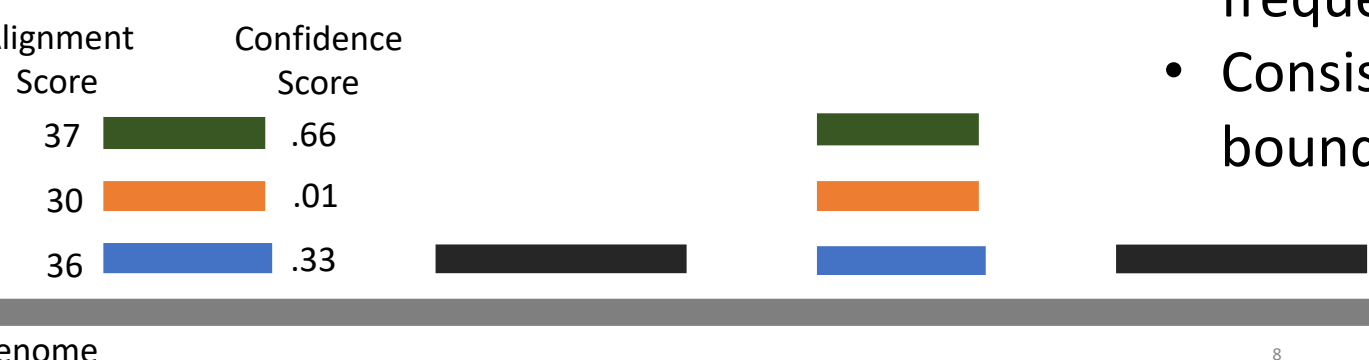
$$score(t, q_i) = \lambda * \log_2 \frac{P(t|q_i)}{P(t|R)} \quad \begin{matrix} \text{[model of homology]} \\ \text{[R=random model]} \end{matrix}$$

So:

$$P(t|q_i) = P(t|R) * 2^{\frac{score(t|q_i)}{\lambda}}$$

$$Conf(q_i|t) = \frac{P(t|R) * 2^{\frac{score(t|q_i)}{\lambda}}}{\sum_j P(t|R) * 2^{\frac{score(t|q_j)}{\lambda}}}$$

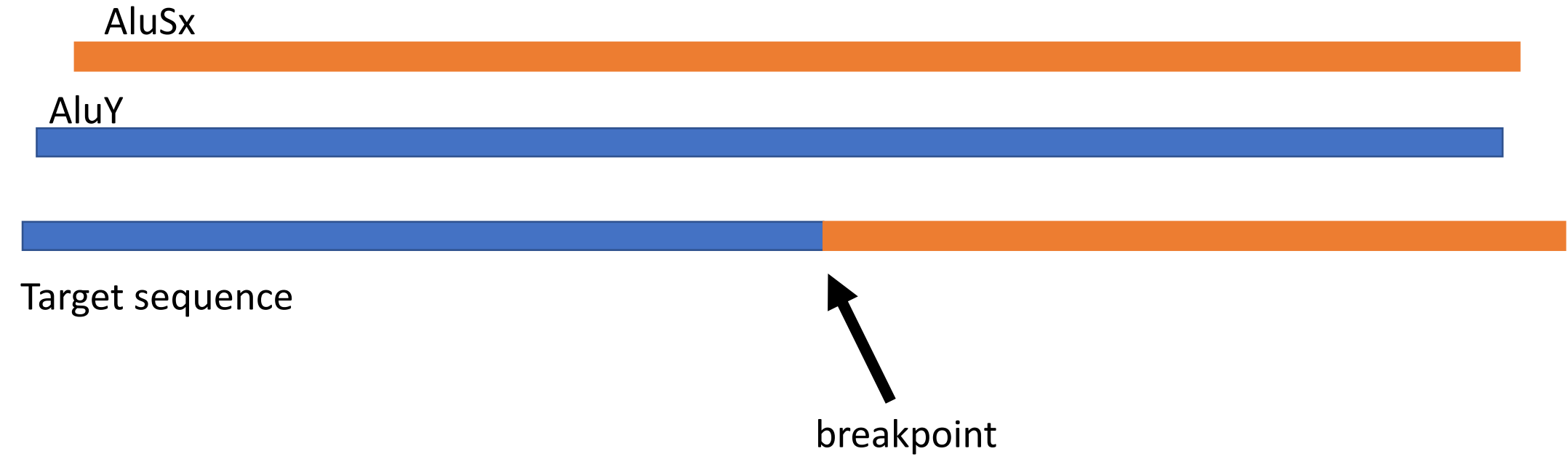
$$Conf(q_i|t) = \frac{2^{\frac{score(t|q_i)}{\lambda}}}{\sum_j 2^{\frac{score(t|q_j)}{\lambda}}}$$



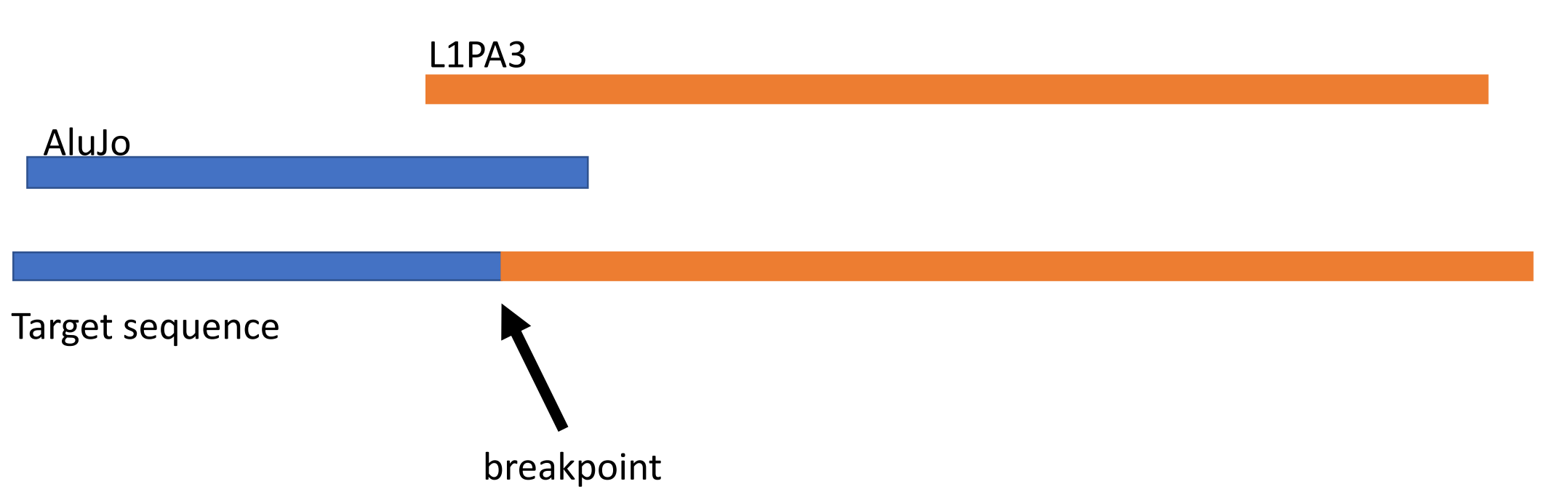
- Assumptions:
- Consistent background frequencies
  - Consistent alignment boundaries

Confidence can also be used to:

1. Identify instances of gene conversion / homologous recombination

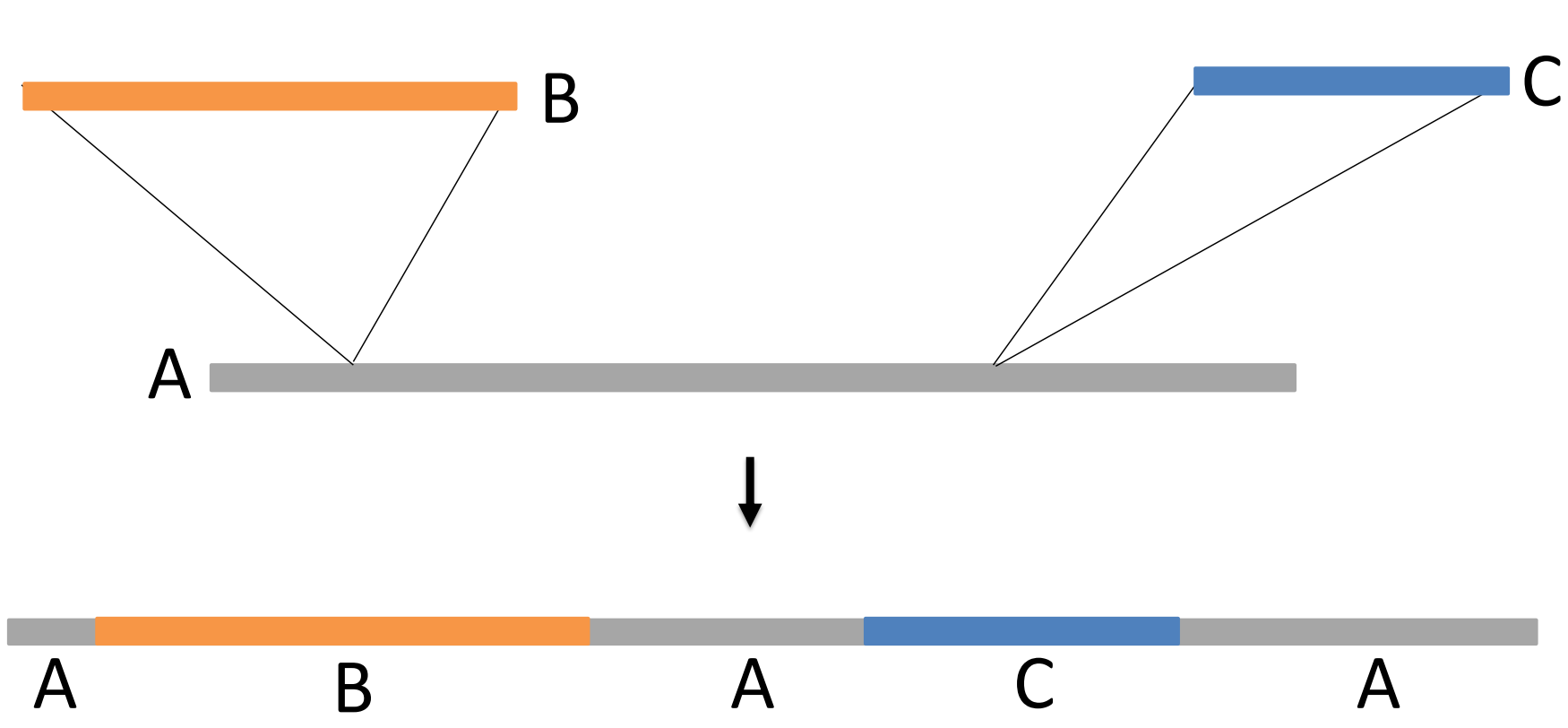


2. Detect boundaries of neighboring partial element matches



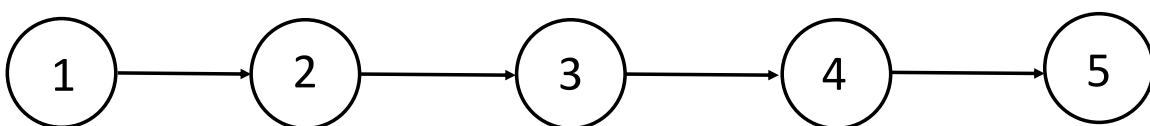
3. Identify inserted elements

- Sequences can also insert themselves inside of other sequences.
- B and C have both inserted themselves inside of A
  - Traditional annotation methods would find 3 occurrences of sequence A
  - We want to find the 3 occurrences of A and recognize them as all belonging to the same original sequence

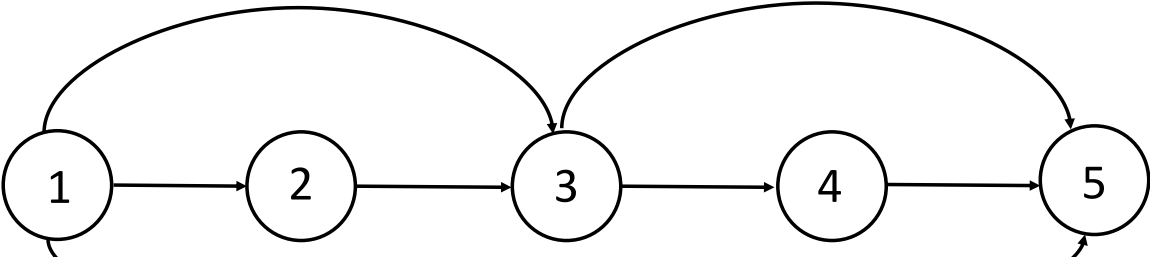


Graph Algorithm Identifies Inserted Elements

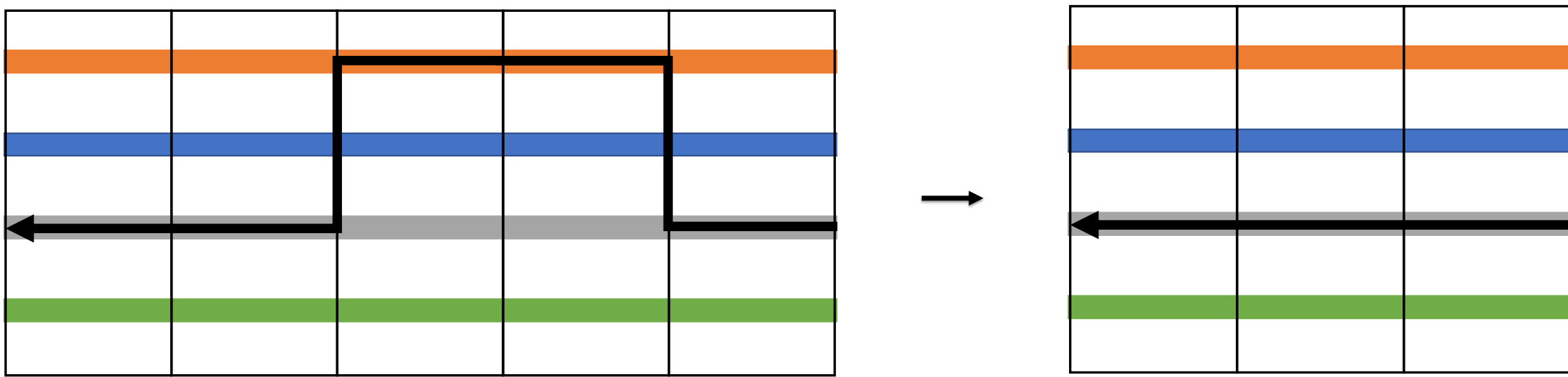
1. Make graph of all sequences labeled in dynamic programming. Each labeled segment becomes a node



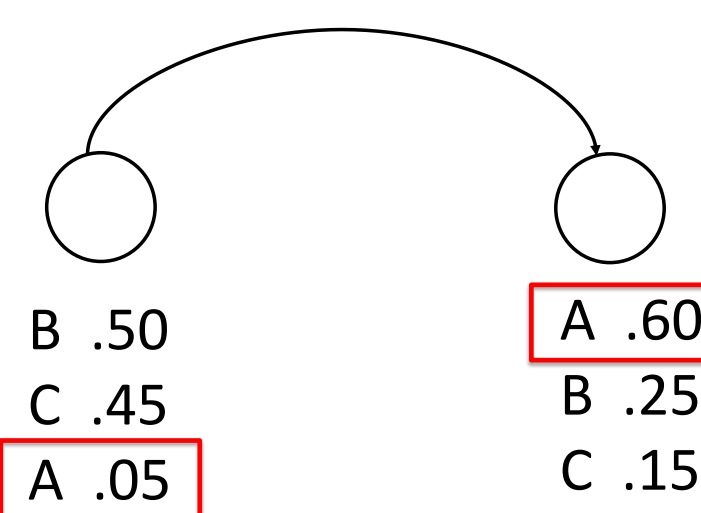
2. Find alternative paths through the graph based on confidence



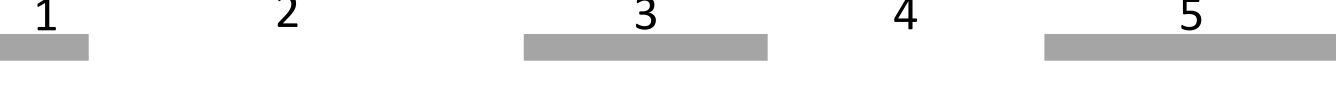
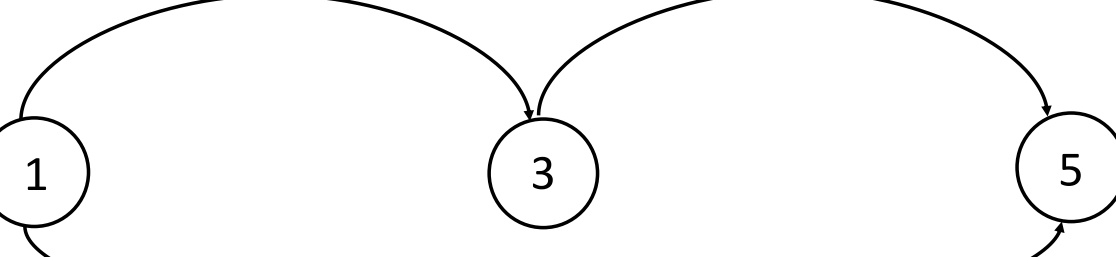
4. Splice corresponding positions out of dynamic programming matrix, stitching the original sequence back together



Edge Creation Rule:  
Create edge if the destination's best label has a confidence > 0.01 in the source

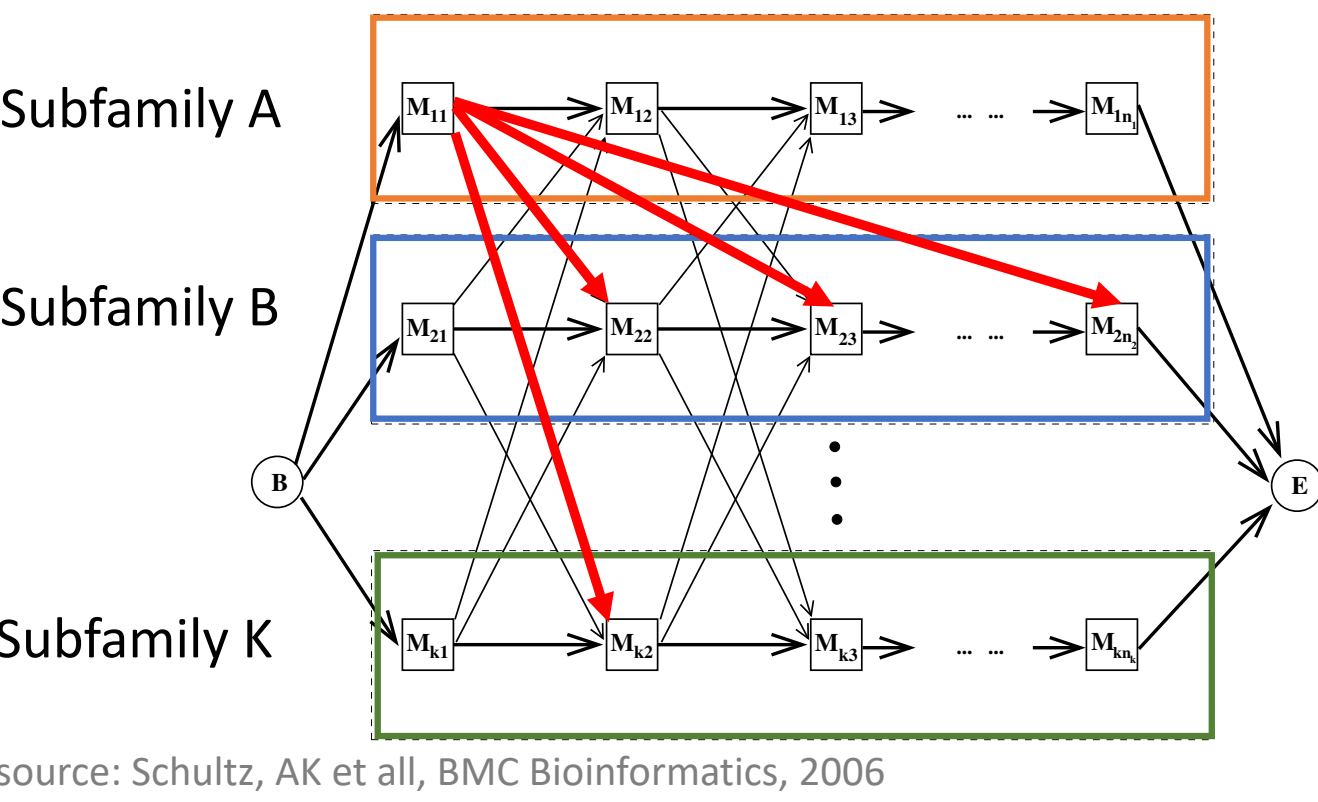


3. Splice out all sequences with no alternate paths

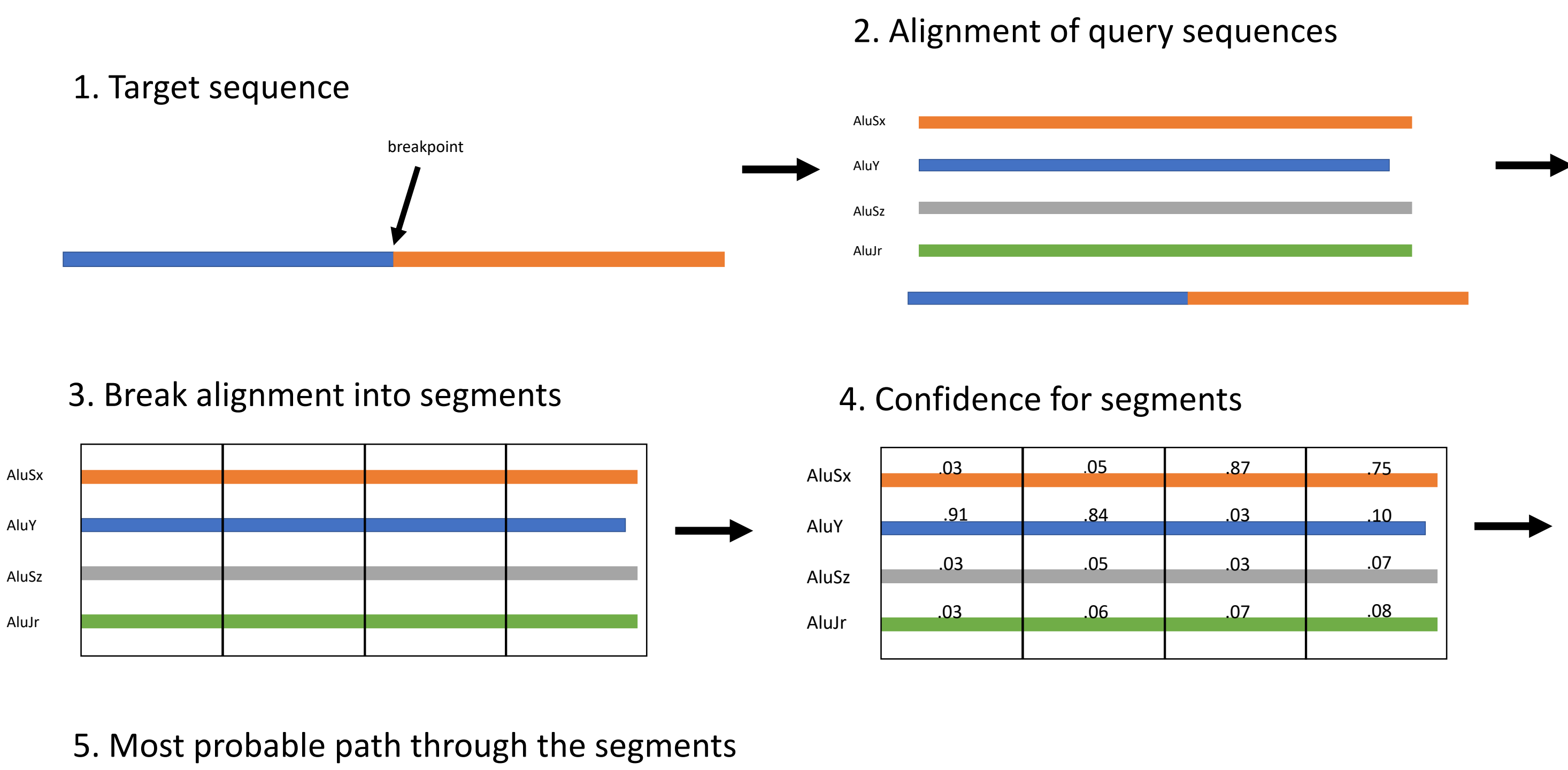


Jumping Profile Hidden Markov Model (ideal)

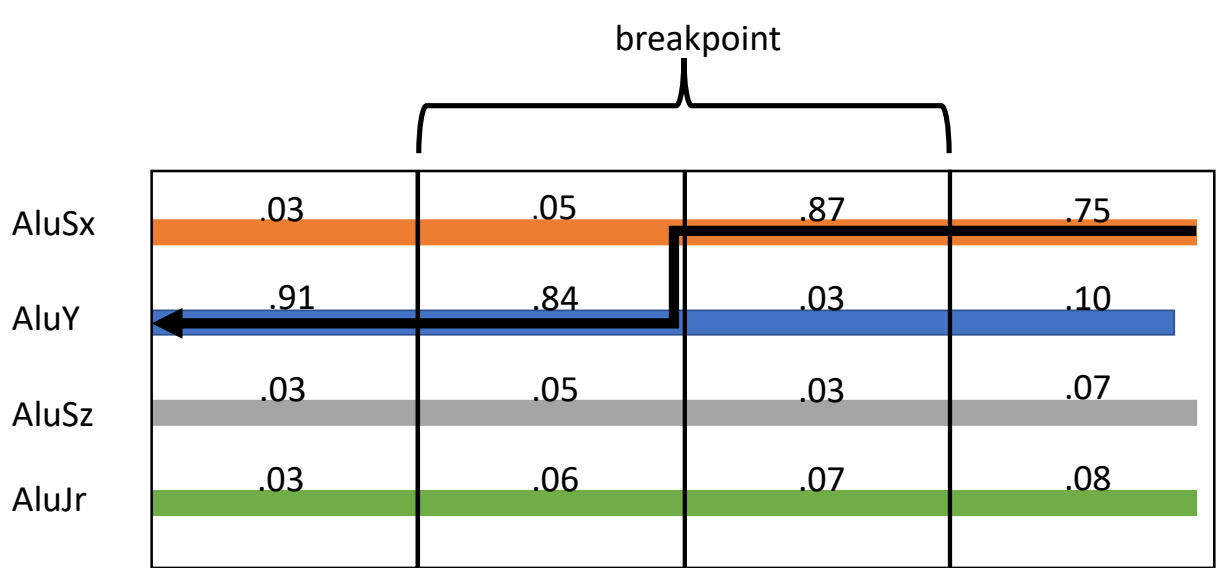
- One approach for identifying breakpoints:
- Sub-model HMM for each subfamily in the alignment
  - Allows transitions (jumps) between submodels
  - Combinatoric explosion of transitions unless the number of possible jumps is constrained



Segmented Confidence Approximates Jumping HMM Approach



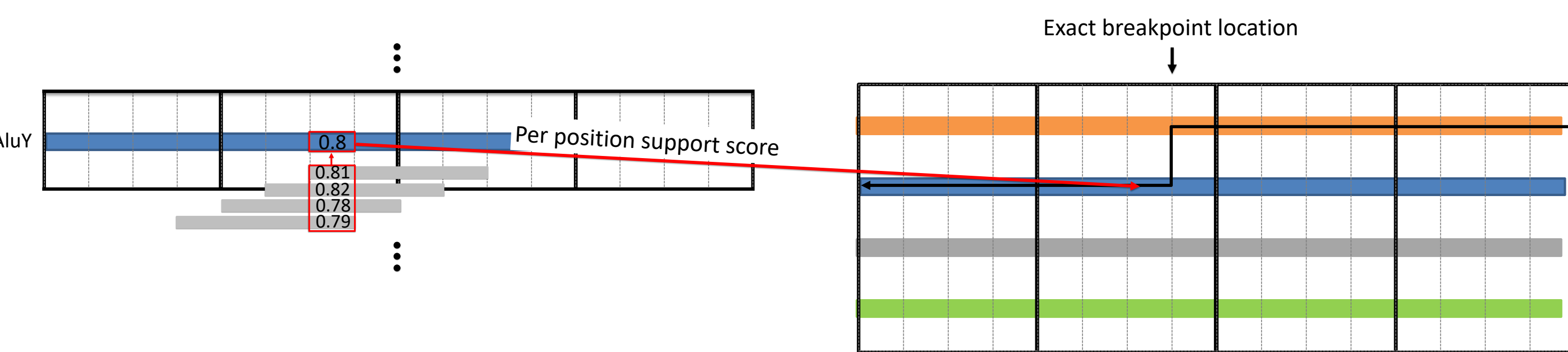
5. Most probable path through the segments



$$S_{i,j} = Conf_{i,j} * \max_j \begin{cases} S_{i-1,j'} * t_m, & \text{if } j \neq j' \\ S_{i-1,j'} * t_s, & \text{if } j = j' \end{cases}$$

Overlapping Confidence Enables Exact Breakpoint Location

Using overlapping segments, starting at every nucleotide, we compute nucleotide specific support scores by averaging all the confidence values for each segment that overlaps the nucleotide position. Per position scores in the dynamic programming matrix, enables us to find the exact breakpoint location.



Example: chr11:11,989,996-11,992,119

