

# Multi-Piece Gap-Affine Penalties for Partial Order Alignment on Graphs

Mathematical Framework and Algorithmic Design

July 3, 2025

## 1 Introduction

Partial order alignment (POA) extends traditional sequence alignment to align sequences against directed acyclic graphs (DAGs) representing multiple sequences. The current POASTA implementation uses standard affine gap penalties, which may over-penalize long insertions and deletions. This document presents a generalized mathematical framework for implementing multi-piece gap-affine penalties that can model complex biological gap length distributions with arbitrary precision.

## 2 Mathematical Framework

### 2.1 Standard Affine Gap Model

The current gap penalty function is defined as:

$$g_{\text{affine}}(k) = \begin{cases} 0 & \text{if } k = 0 \\ \alpha + \beta \cdot (k - 1) & \text{if } k > 0 \end{cases} \quad (1)$$

where  $\alpha$  is the gap opening penalty,  $\beta$  is the gap extension penalty, and  $k$  is the gap length.

### 2.2 Generalized Multi-Piece Gap-Affine Model

The generalized n-piece gap-affine model uses multiple breakpoints  $\mathbf{k} = [k_1, k_2, \dots, k_{n-1}]$  to create piecewise linear gap penalties:

$$g_{\text{multi-piece}}(\ell) = \begin{cases} 0 & \text{if } \ell = 0 \\ \alpha + \sum_{i=1}^{j-1} \beta_i \cdot (k_i - k_{i-1}) + \beta_j \cdot (\ell - k_{j-1}) & \text{if } k_{j-1} < \ell \leq k_j \end{cases} \quad (2)$$

where:

$$\alpha = \text{gap opening penalty} \quad (3)$$

$$\mathbf{k} = [k_0, k_1, k_2, \dots, k_{n-1}, k_n] \text{ with } k_0 = 1, k_n = \infty \quad (4)$$

$$\beta = [\beta_1, \beta_2, \dots, \beta_n] \text{ extension penalties for each piece} \quad (5)$$

$$j = \text{piece index such that } k_{j-1} < \ell \leq k_j \quad (6)$$

### Special Cases:

$$\text{Standard affine (1-piece): } \mathbf{k} = [1, \infty], \boldsymbol{\beta} = [\beta] \quad (7)$$

$$\text{Two-piece: } \mathbf{k} = [1, k_0, \infty], \boldsymbol{\beta} = [\beta_1, \beta_2] \quad (8)$$

$$\text{Three-piece: } \mathbf{k} = [1, k_0, k_1, \infty], \boldsymbol{\beta} = [\beta_1, \beta_2, \beta_3] \quad (9)$$

Typically,  $\beta_1 > \beta_2 > \beta_3 > \dots$  to progressively reduce penalty for longer gaps.

## 2.3 Biological Justification

Empirical studies of protein structural alignments reveal complex gap length distributions that can be modeled with multiple linear segments:

$$\log P(\text{gap length} = \ell) = \begin{cases} a_1 + b_1 \cdot \ell & \text{if } 1 \leq \ell \leq 3 \\ a_2 + b_2 \cdot \ell & \text{if } 4 \leq \ell \leq 20 \\ a_3 + b_3 \cdot \ell & \text{if } \ell > 20 \end{cases} \quad (10)$$

where  $|b_1| > |b_2| > |b_3|$ , reflecting:

- Short gaps (1-3): High penalty due to structural constraints
- Medium gaps (4-20): Moderate penalty for common evolutionary events
- Long gaps ( $>20$ ): Low penalty for large structural rearrangements

## 3 Graph Alignment Formulation

### 3.1 Problem Definition

Let  $G = (V, E)$  be a directed acyclic graph representing a partial order alignment, and let  $q = q_1 q_2 \dots q_n$  be a query sequence. The goal is to find an optimal alignment path through  $G$  that aligns  $q$  to the graph.

### 3.2 State Space

Define the generalized state space as:

$$\mathcal{S} = \{(i, v, s, j, \ell) : i \in [0, n], v \in V, s \in \{M, I, D\}, j \in [1, p], \ell \in \mathbb{N}\} \quad (11)$$

where:

$$i = \text{position in query sequence} \quad (12)$$

$$v = \text{current node in graph} \quad (13)$$

$$s = \text{alignment state} \in \{M, I, D\} \quad (14)$$

$$j = \text{current piece index} \in [1, p] \quad (15)$$

$$\ell = \text{gap length within current piece} \quad (16)$$

The alignment states are:

$$M = \text{match/mismatch state} \quad (17)$$

$$I = \text{insertion state (gap in graph)} \quad (18)$$

$$D = \text{deletion state (gap in query)} \quad (19)$$

### 3.3 Dynamic Programming Recurrence Relations

Let  $\text{OPT}(i, v, s, j, \ell)$  denote the optimal alignment score for state  $(i, v, s, j, \ell)$ .

#### 3.3.1 Match State

$$\text{OPT}(i, v, \text{M}) = \min_{u \in \text{pred}(v)} \{ \text{OPT}(i-1, u, \text{M}) + \sigma(q_i, v), \quad (20)$$

$$\min_{j, \ell} \text{OPT}(i-1, u, \text{I}, j, \ell) + \sigma(q_i, v), \quad (21)$$

$$\min_{j, \ell} \text{OPT}(i-1, u, \text{D}, j, \ell) + \sigma(q_i, v) \} \quad (22)$$

where  $\sigma(q_i, v)$  is the substitution score for aligning query character  $q_i$  with graph node  $v$ .

#### 3.3.2 Insertion States

**Starting new gap:**

$$\text{OPT}(i, v, \text{I}, 1, 1) = \text{OPT}(i-1, v, \text{M}) + \alpha + \beta_1 \quad (23)$$

**Extending within same piece:**

$$\text{OPT}(i, v, \text{I}, j, \ell+1) = \text{OPT}(i-1, v, \text{I}, j, \ell) + \beta_j \quad \text{if } \ell+1 \leq k_j - k_{j-1} \quad (24)$$

**Transitioning to next piece:**

$$\text{OPT}(i, v, \text{I}, j+1, 1) = \text{OPT}(i-1, v, \text{I}, j, k_j - k_{j-1}) + \beta_{j+1} \quad \text{if } j < p \quad (25)$$

#### 3.3.3 Deletion States

**Starting new gap:**

$$\text{OPT}(i, v, \text{D}, 1, 1) = \min_{u \in \text{pred}(v)} \text{OPT}(i, u, \text{M}) + \alpha + \beta_1 \quad (26)$$

**Extending within same piece:**

$$\text{OPT}(i, v, \text{D}, j, \ell+1) = \min_{u \in \text{pred}(v)} \text{OPT}(i, u, \text{D}, j, \ell) + \beta_j \quad \text{if } \ell+1 \leq k_j - k_{j-1} \quad (27)$$

**Transitioning to next piece:**

$$\text{OPT}(i, v, \text{D}, j+1, 1) = \min_{u \in \text{pred}(v)} \text{OPT}(i, u, \text{D}, j, k_j - k_{j-1}) + \beta_{j+1} \quad \text{if } j < p \quad (28)$$

## 4 A\* Algorithm Adaptation

### 4.1 Heuristic Function

The admissible heuristic function for the remaining alignment cost is:

$$h(i, v, s, j, \ell) = h_{\text{gap}}(n-i) + h_{\text{sub}}(i, v) \quad (29)$$

where  $h_{\text{gap}}(r)$  is the minimum gap cost for  $r$  remaining characters using the multi-piece gap cost function:

$$h_{\text{gap}}(r) = g_{\text{multi-piece}}(r) \quad (30)$$

and  $h_{\text{sub}}(i, v)$  is the minimum substitution cost from node  $v$  to the end of the graph.

## 4.2 Priority Function

The A\* priority function is:

$$f(i, v, s, j, \ell) = g(i, v, s, j, \ell) + h(i, v, s, j, \ell) \quad (31)$$

where  $g(i, v, s, j, \ell)$  is the actual cost to reach state  $(i, v, s, j, \ell)$ .

## 5 Complexity Analysis

### 5.1 Time Complexity

The time complexity of the multi-piece gap-affine algorithm is:

$$O(|V| \cdot |E| \cdot n \cdot k_{\max} \cdot p) \quad (32)$$

where  $k_{\max}$  is the maximum gap length considered and  $p$  is the number of pieces. This represents a factor of  $k_{\max} \cdot p$  increase over the standard affine model.

### 5.2 Space Complexity

The space complexity is:

$$O(|V| \cdot n \cdot k_{\max} \cdot p) \quad (33)$$

The additional space requirement can be managed through:

- Limiting  $k_{\max}$  to biologically reasonable values (e.g., 20-50)
- Using sparse representations for gap length states
- Implementing state pruning strategies
- Limiting  $p$  to 2-4 pieces for practical applications

## 6 Implementation Considerations

### 6.1 Parameter Selection

Based on empirical studies, recommended parameter configurations are:

**Two-piece model:**

$$\mathbf{k} = [1, 3, \infty] \quad (34)$$

$$\boldsymbol{\beta} = [2, 1] \quad (35)$$

$$\alpha = 6 \quad (36)$$

**Three-piece model:**

$$\mathbf{k} = [1, 3, 20, \infty] \quad (37)$$

$$\boldsymbol{\beta} = [2, 1, 0] \quad (38)$$

$$\alpha = 6 \quad (39)$$

## 6.2 Memory Optimization

To manage memory requirements:

1. Use blocked storage with separate blocks for different gap length ranges
2. Implement lazy state creation - only create states when needed
3. Employ sparse representations for piece-length combinations
4. Use bit-packing for small gap lengths within pieces
5. Implement dynamic piece pruning based on gap length distributions

## 6.3 Algorithmic Optimizations

1. **Piece Transition Caching:** Pre-compute transition costs between pieces
2. **Gap Length Pruning:** Dynamically limit maximum gap length based on heuristics
3. **State Aggregation:** Combine similar gap states to reduce search space
4. **Adaptive Piece Selection:** Choose number of pieces based on input characteristics

## 7 Conclusion

The generalized multi-piece gap-affine penalty model provides a flexible and biologically realistic approach to gap scoring in partial order alignment. While it increases computational complexity by a factor of  $k_{\max} \cdot p$ , the linear scaling with the number of pieces makes it practical for small  $p$  (2-4 pieces) with careful implementation.

Key advantages of the generalized framework:

- **Unified Design:** All gap penalty models handled by the same algorithm
- **Biological Realism:** Can model complex empirical gap distributions
- **Extensibility:** Easy to add new pieces without algorithmic changes
- **Backward Compatibility:** Standard affine gaps are a special case
- **Research Flexibility:** Enables experimentation with different penalty schemes

The mathematical framework presented here provides a solid foundation for implementing this generalized model in the POASTA system, with clear algorithmic specifications and complexity analysis to guide the implementation process. The unified approach eliminates the need for separate implementations of 1-piece, 2-piece, and 3-piece models, resulting in cleaner, more maintainable code.