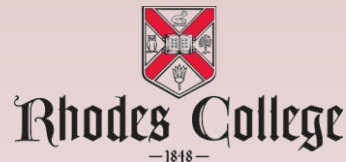# COMP 355
# Advanced Algorithms

**Dynamic Programming:**
**Space Efficient Sequence Alignment**
**Section 6.7(KT)**

# Edit Distance

Applications.

- Basis for Unix diff.

- Speech recognition.

- Computational biology.

Edit distance.  [Levenshtein 1966, Needleman-Wunsch 1970]

- Gap penalty $\delta$; mismatch penalty $\alpha_{pq}$.

- Cost = sum of gap and mismatch penalties.



$$\alpha_{TC} + \alpha_{GT} + \alpha_{AG} + 2\alpha_{CA}$$

$$2\delta + \alpha_{CA}$$

# Sequence Alignment: Problem Structure

Def.  OPT($i$, $j$) = min cost of aligning strings $x_1 x_2 \ldots x_i$ and $y_1 y_2 \ldots y_j$.

- Case 1:  OPT matches $x_i$-$y_j$.
  - pay mismatch for $x_i$-$y_j$  + min cost of aligning two strings $x_1 x_2 \ldots x_{i-1}$ and $y_1 y_2 \ldots y_{j-1}$
- Case 2a:  OPT leaves $x_i$ unmatched.
  - pay gap for $x_i$ and min cost of aligning $x_1 x_2 \ldots x_{i-1}$ and $y_1 y_2 \ldots y_j$
- Case 2b:  OPT leaves $y_j$ unmatched.
  - pay gap for $y_j$ and min cost of aligning $x_1 x_2 \ldots x_i$ and $y_1 y_2 \ldots y_{j-1}$

$$OPT(i,\ j) = \begin{cases} j\delta & \text{if } i = 0 \\ \min \begin{cases} \alpha_{x_i y_j} + OPT(i-1,\ j-1) \\ \delta + OPT(i-1,\ j) \\ \delta + OPT(i,\ j-1) \end{cases} & \text{otherwise} \\ i\delta & \text{if } j = 0 \end{cases}$$

# Sequence Alignment:  Algorithm

```
Sequence-Alignment(m, n, x₁x₂...xₘ, y₁y₂...yₙ, δ, α) {
    for i = 0 to m
        M[0, i] = iδ
    for j = 0 to n
        M[j, 0] = jδ

    for i = 1 to m
        for j = 1 to n
            M[i, j] = min(α[xᵢ, yⱼ] + M[i-1, j-1],
                          δ + M[i-1, j],
                          δ + M[i, j-1])
    return M[m, n]
}
```

- Analysis.  $\Theta(mn)$ time and space.
- English words or sentences:  m, n $\leq$ 10.
- Computational biology:  m = n = 100,000. 10 billions ops OK, but 10GB array?

# Sequence Alignment:  Linear Space

Q.  Can we avoid using quadratic space?

Easy.  Optimal value in O(m + n) space and O(mn) time.

- Compute OPT(i, •) from OPT(i-1, •).
- No longer a simple way to recover alignment itself.
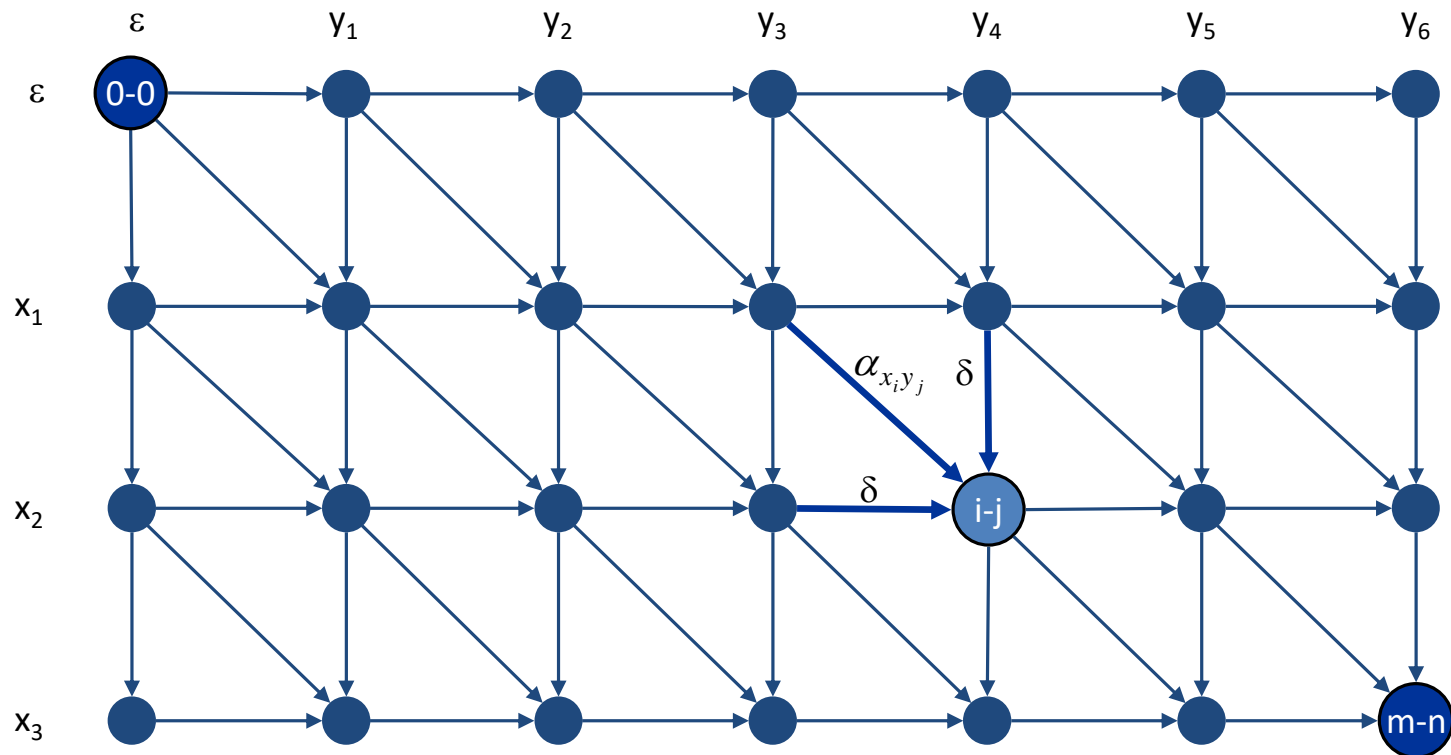
Theorem.  [Hirschberg 1975] Optimal alignment in $O(m + n)$ space and $O(mn)$ time.

- Clever combination of divide-and-conquer and dynamic programming.
- Inspired by idea of Savitch from complexity theory.

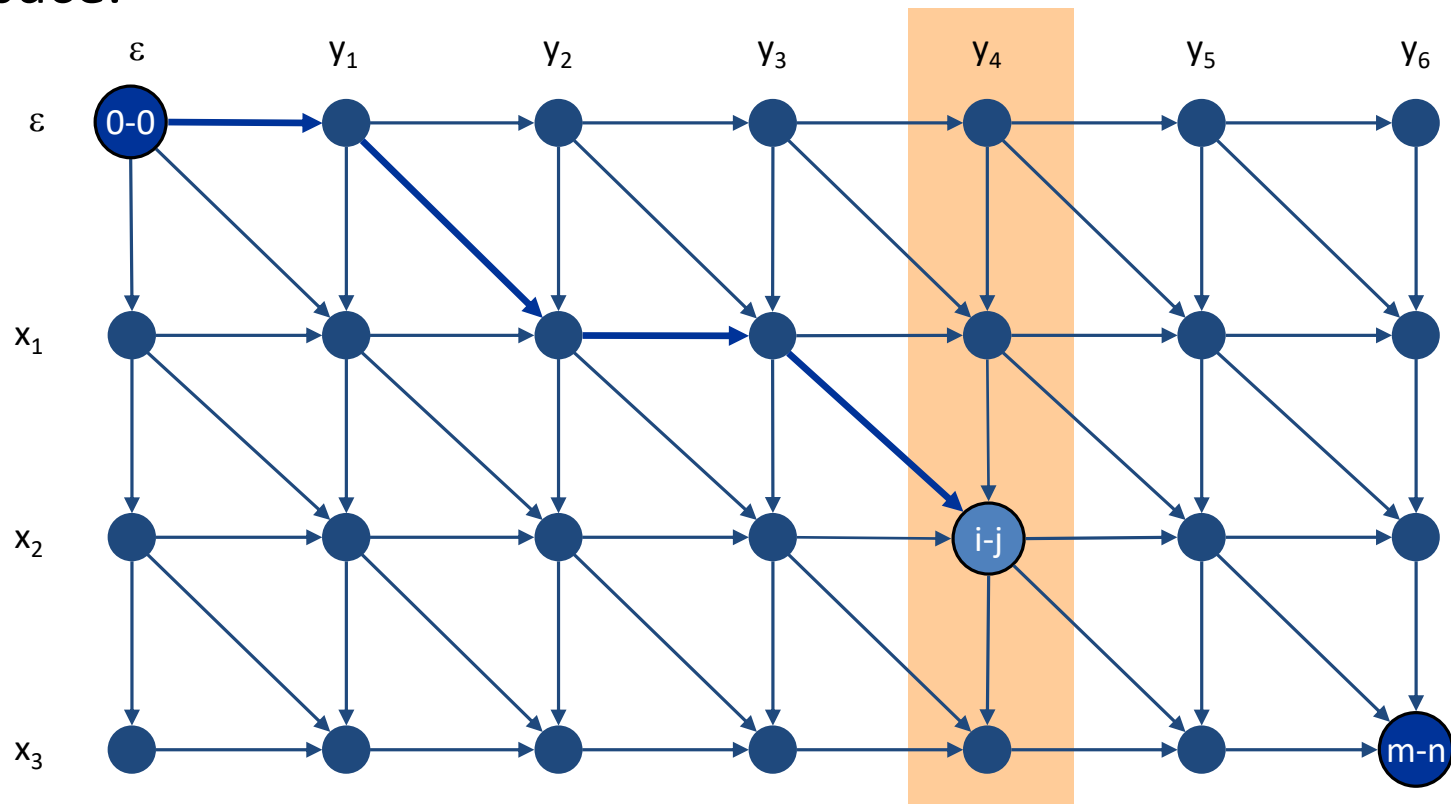# Sequence Alignment:  Linear Space

Edit distance graph.

- Let f(i, j) be shortest path from (0,0) to (i, j).

- Observation:  f(i, j) = OPT(i, j).

# Sequence Alignment:  Linear Space

Edit distance graph.

- Let $f(i, j)$ be shortest path from $(0,0)$ to $(i, j)$.

- Can compute $f(\bullet, j)$ for any $j$ in $O(mn)$ time and $O(m + n)$ space.

# Space-Efficient Alignment Algorithm

```
Space-Efficient-Alignment(X,Y)
    Array B[0 ... m, 0 ... 1]
    Initialize B[i,0] = i δ (just as in column 0 of A)
    For j = 1, ..., n
        B[0, 1] = jδ (since this corresponds to entry A[0, j])
        For i = 1, ..., m
            B[i, 1] = min[α_{x_i y_j} + B[i − 1, 0],
                          δ + B[i − 1, 1], δ + B[i, 0]].
        Endfor
        Move column 1 of B to column 0 to make room for next iteration:
            Update B[i, 0] = B[i, 1] for each i
    Endfor
```
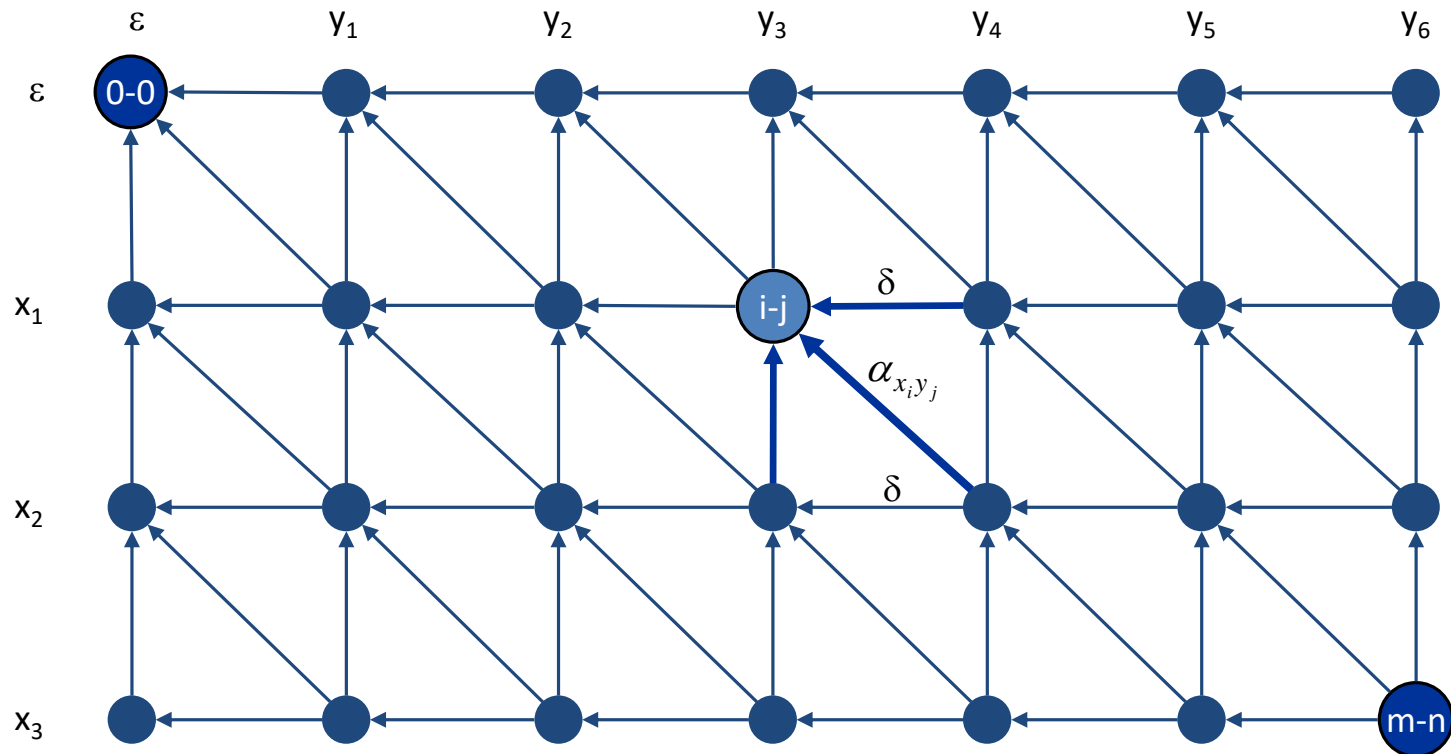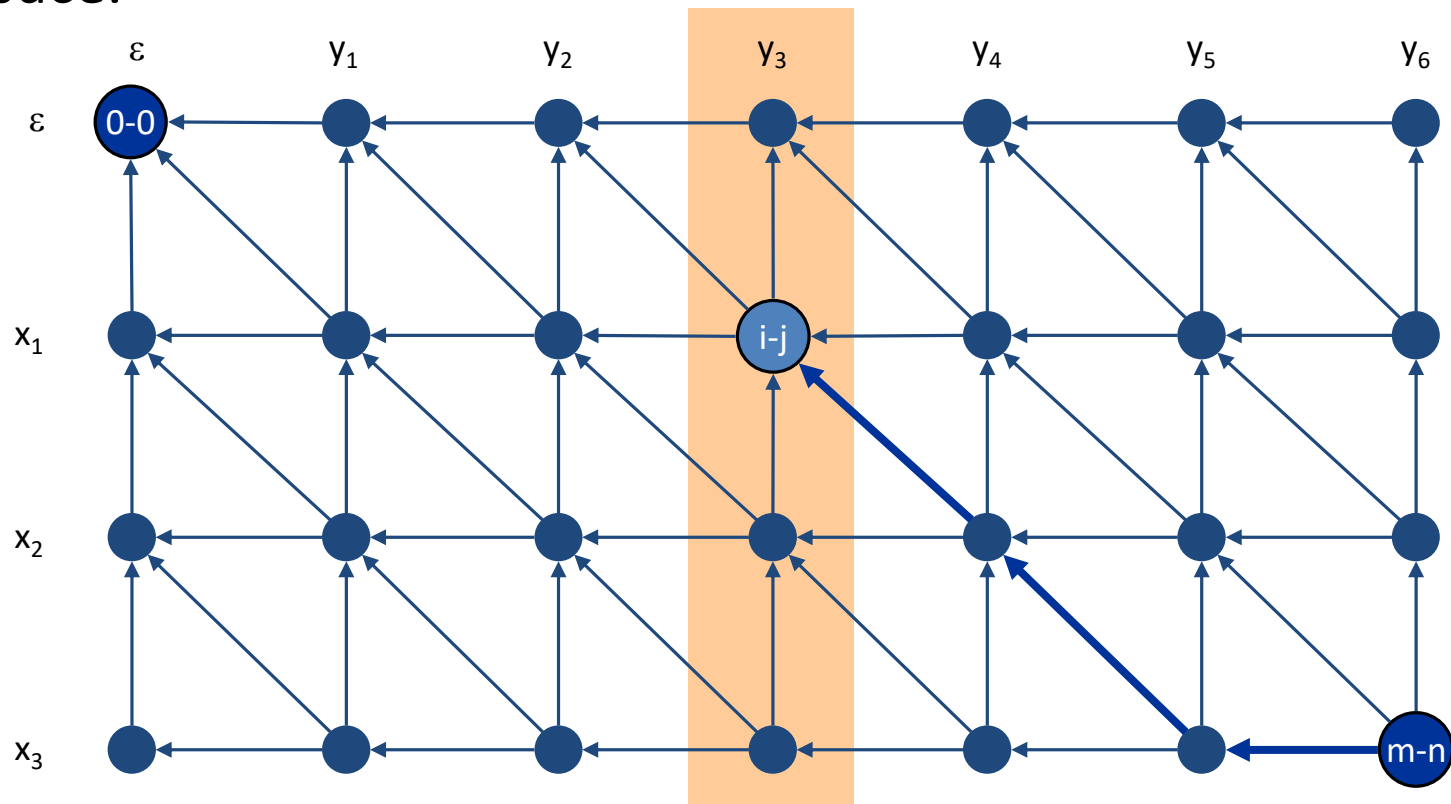
# Sequence Alignment:  Linear Space

Edit distance graph.

- Let g(i, j) be shortest path from (i, j) to (m, n).

- Can compute by reversing the edge orientations and inverting the roles of (0, 0) and (m, n)
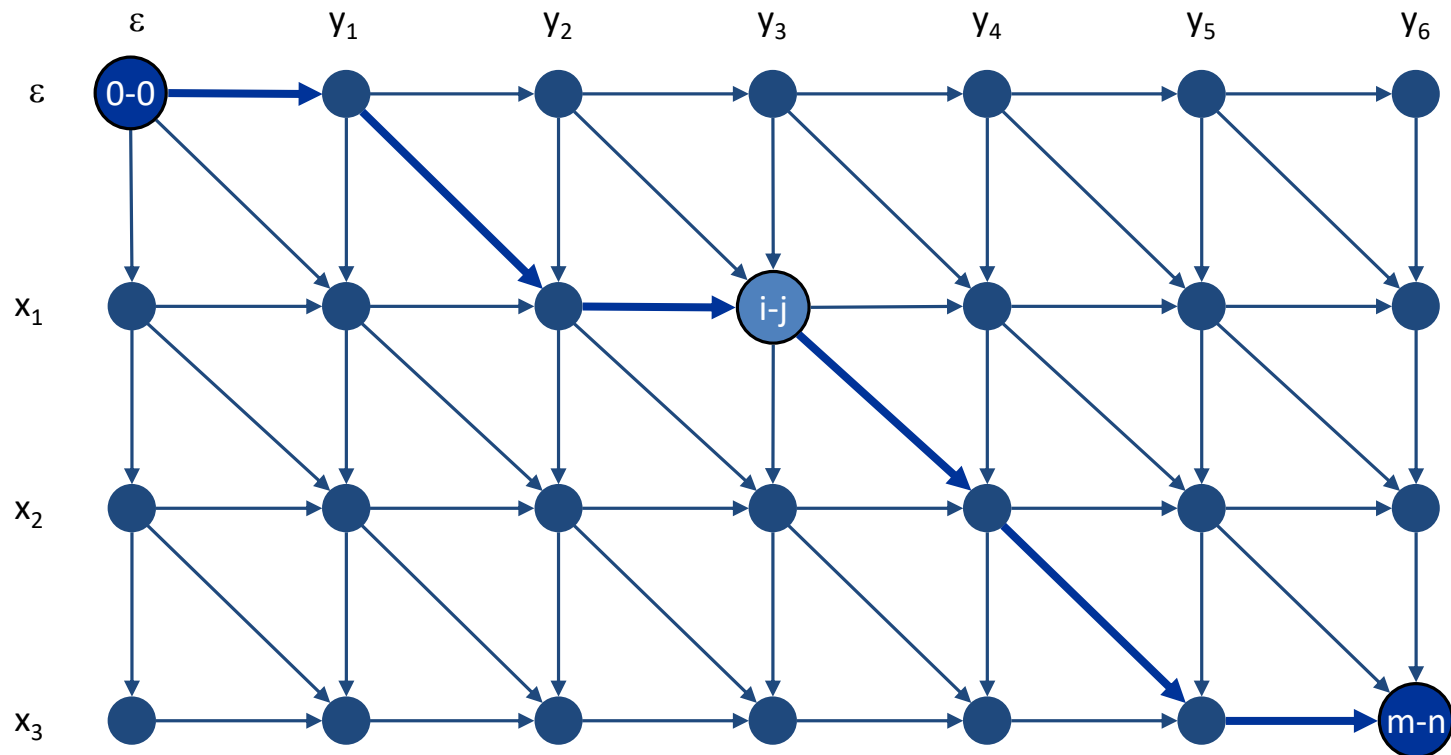
# Sequence Alignment: Linear Space

Edit distance graph.

- Let g(i, j) be shortest path from (i, j) to (m, n).
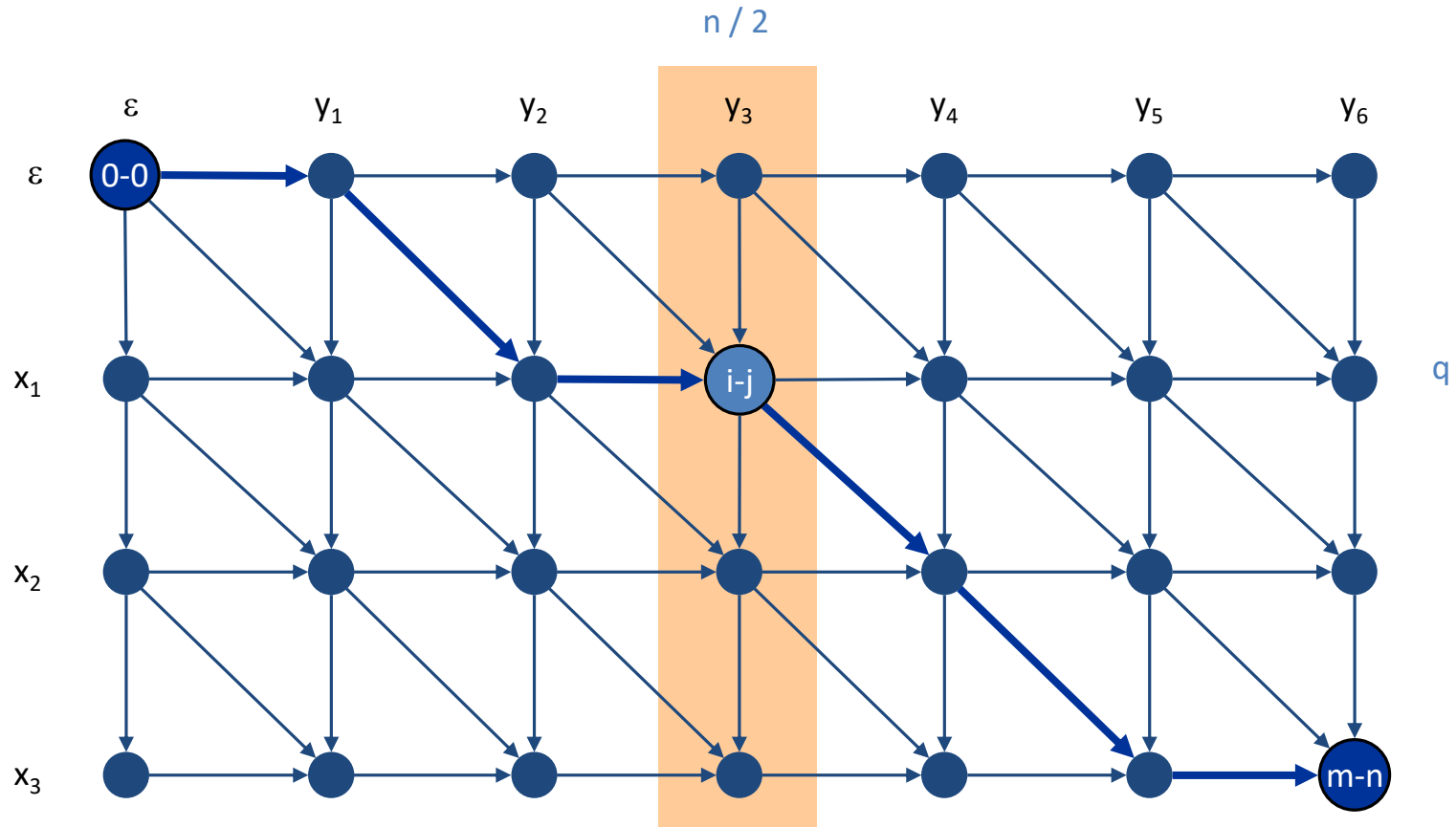- Can compute g(•, j) for any j in O(mn) time and O(m + n) space.

# Sequence Alignment:  Linear Space

Observation 1.  The cost of the shortest path that uses (i, j) is f(i, j) + g(i, j).

Observation 2.  let q be an index that minimizes $f(q, n/2) + g(q, n/2)$. Then, the shortest path from $(0, 0)$ to $(m, n)$ uses $(q, n/2)$.
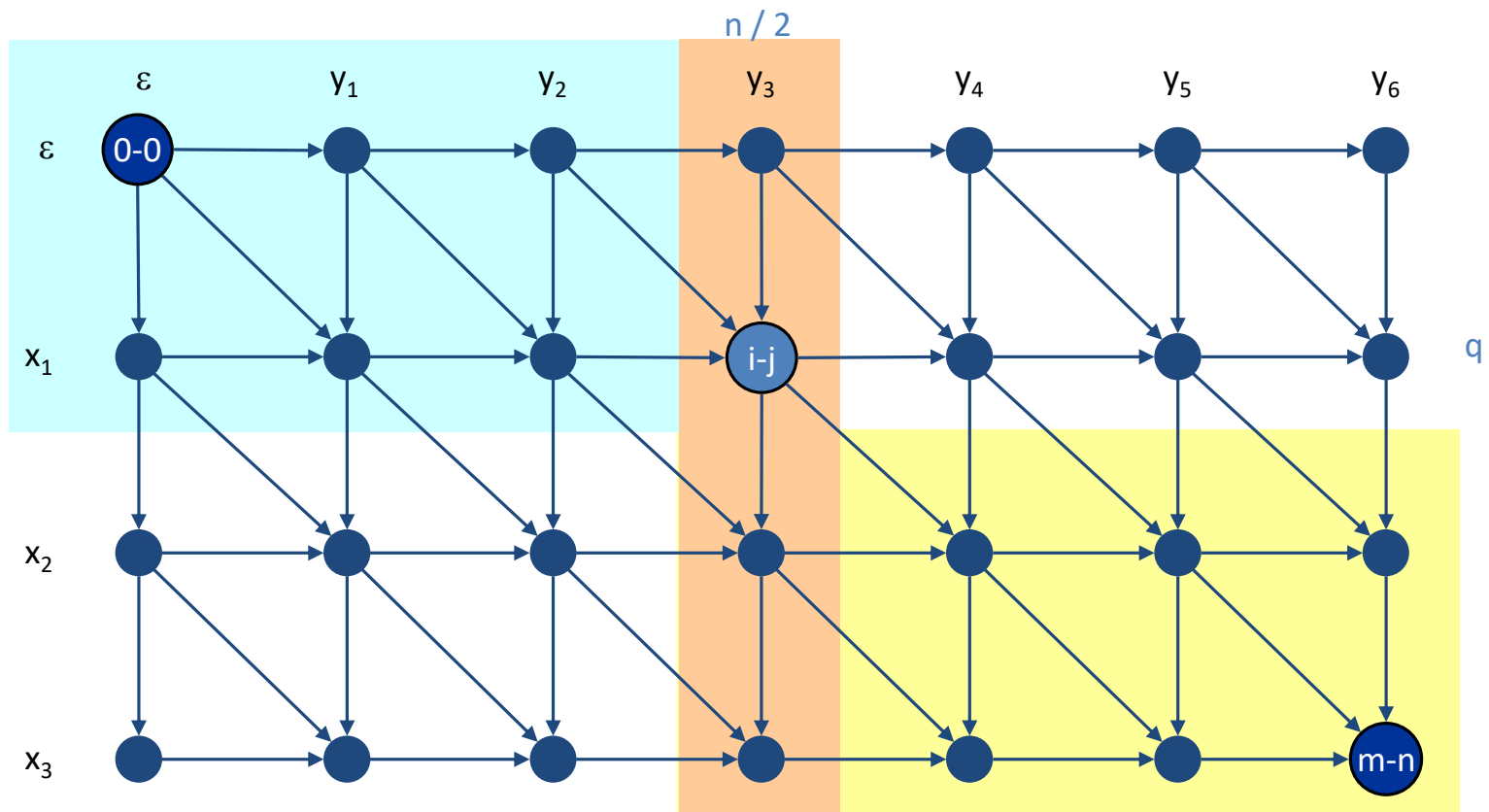
# Sequence Alignment:  Linear Space

Divide:  find index q that minimizes $f(q, n/2) + g(q, n/2)$ using DP.

- Align $x_q$ and $y_{n/2}$.

Conquer:  recursively compute optimal alignment in each piece.

# Sequence Alignment:
# Running Time Analysis Warmup

**Theorem.** Let $T(m, n)$ = max running time of algorithm on strings of length at most $m$ and $n$. $T(m, n) = O(mn \log n)$.

$$T(m, n) \leq 2T(m, n/2) + O(mn) \implies T(m, n) = O(mn \log n)$$

**Remark.** Analysis is not tight because two sub-problems are of size (q, n/2) and (m - q, n/2). In next slide, we save $\log n$ factor.

# Sequence Alignment: Running Time Analysis

Theorem.  Let $T(m,n)$ = max running time of algorithm on strings of length m and n. $T(m,n) = O(mn)$.

Pf.  (by induction on n)
- $O(mn)$ time to compute $f(\bullet, n/2)$ and $g(\bullet, n/2)$ and find index $q$.
- $T(q, n/2) + T(m - q, n/$
- Choose constant $c$ so that:

$$
\begin{aligned}
T(m,\ 2) &\leq cm \\
T(2,\ n) &\leq cn \\
T(m,\ n) &\leq cmn + T(q,\ n/2) + T(m-q,\ n/2)
\end{aligned}
$$

- Basis cases: $m = 2$ or $n = 2$.
- Inductive hypothesis:  $T(m, n)$

$$
\begin{aligned}
T(m,n) &\leq T(q,n/2) + T(m-q,n/2) + cmn \\
&\leq 2cqn/2 + 2c(m-q)n/2 + cmn \\
&= cqn + cmn - cqn + cmn \\
&= 2cmn
\end{aligned}
$$

# Divide-and-Conquer Alignment Algorithm

```
Divide-and-Conquer-Alignment(X,Y)
  Let m be the number of symbols in X
  Let n be the number of symbols in Y
  If m ≤ 2 or n ≤ 2 then
      Compute optimal alignment using Alignment(X,Y)
  Call Space-Efficient-Alignment(X,Y[1:n/2]),
      obtaining array B
  Call Backward-Space-Efficient-Alignment(X,Y[n/2+1:n]),
      obtaining array B'
  Let q be the index minimizing B[q,1] + B'[q,n]
  Add (q,n/2) to global list P
  Divide-and-Conquer-Alignment(X[1:q],Y[1:n/2])
  Divide-and-Conquer-Alignment(X[q+1:n],Y[n/2+1:n])
  Return P
```