

# Soft-DTW: A Differentiable Loss Function for Time Series

Marco Cuturi



Mathieu Blondel



In proceedings of ICML'17

Follow pollution levels  
in real time in your city



From: Plume App

**Ground truth  
(reality)**



From: Plume App

**Ground truth  
(reality)**



**How wrong  
was this  
prediction?**

**This depends  
on the loss  
function used  
to train the  
algorithm.**

From: Plume App

- In this talk we propose to use the celebrated **Dynamic Time Warping** discrepancy as a loss.
- Loss functions should be **differentiable**. We show that an **appropriate smoothing** , **soft-DTW**, helps
- We apply this to **several problems**:
  - Computation of **barycenters**,
  - Clustering of time series,
  - Learning with structured (time series) output

Groups  
of time series  
(reality)

Follow pollution levels  
in real time in your city

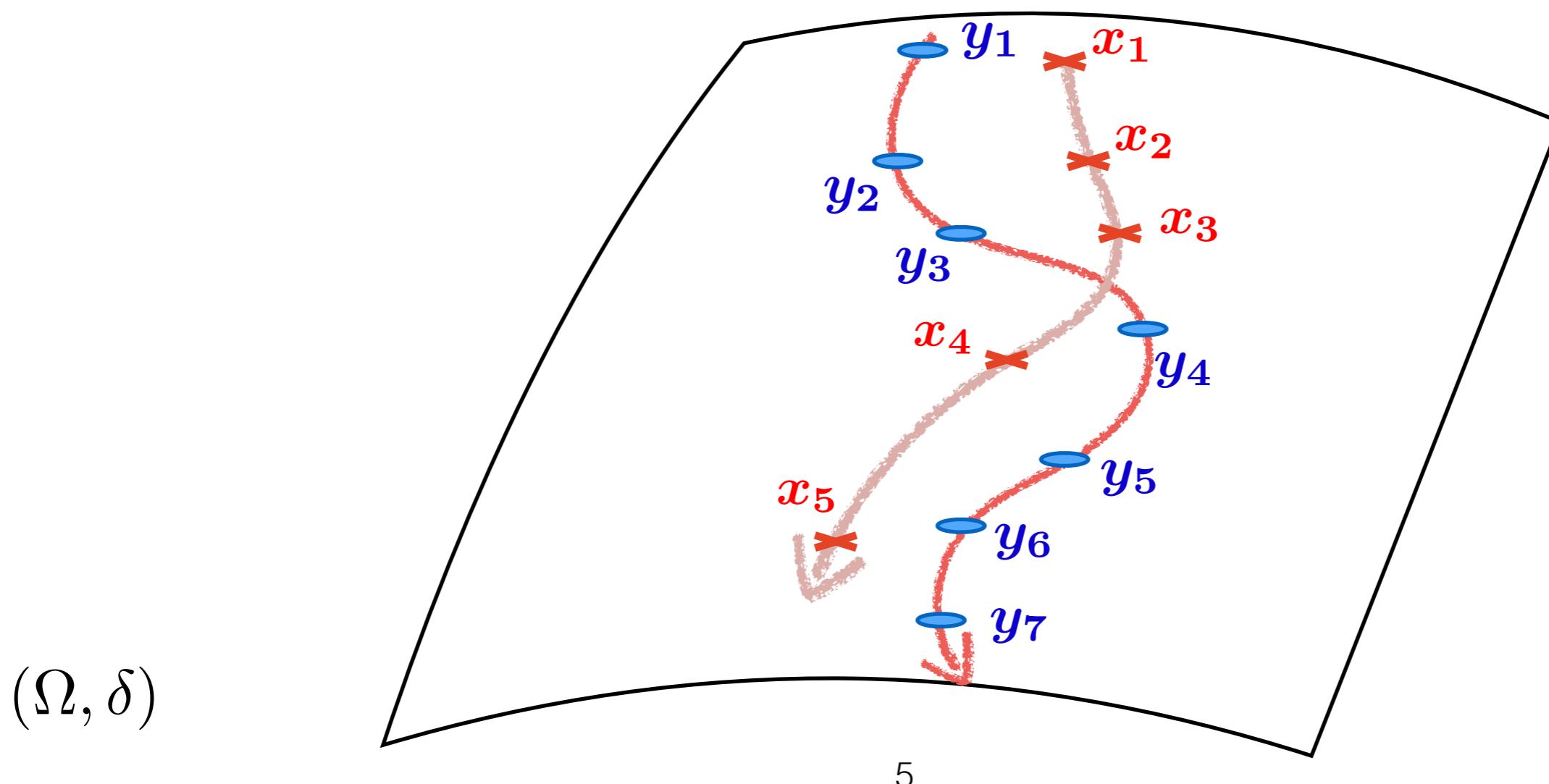
# 0. The DTW Geometry

## 1. Soft-DTW

## 2. Soft-DTW as a Loss Function

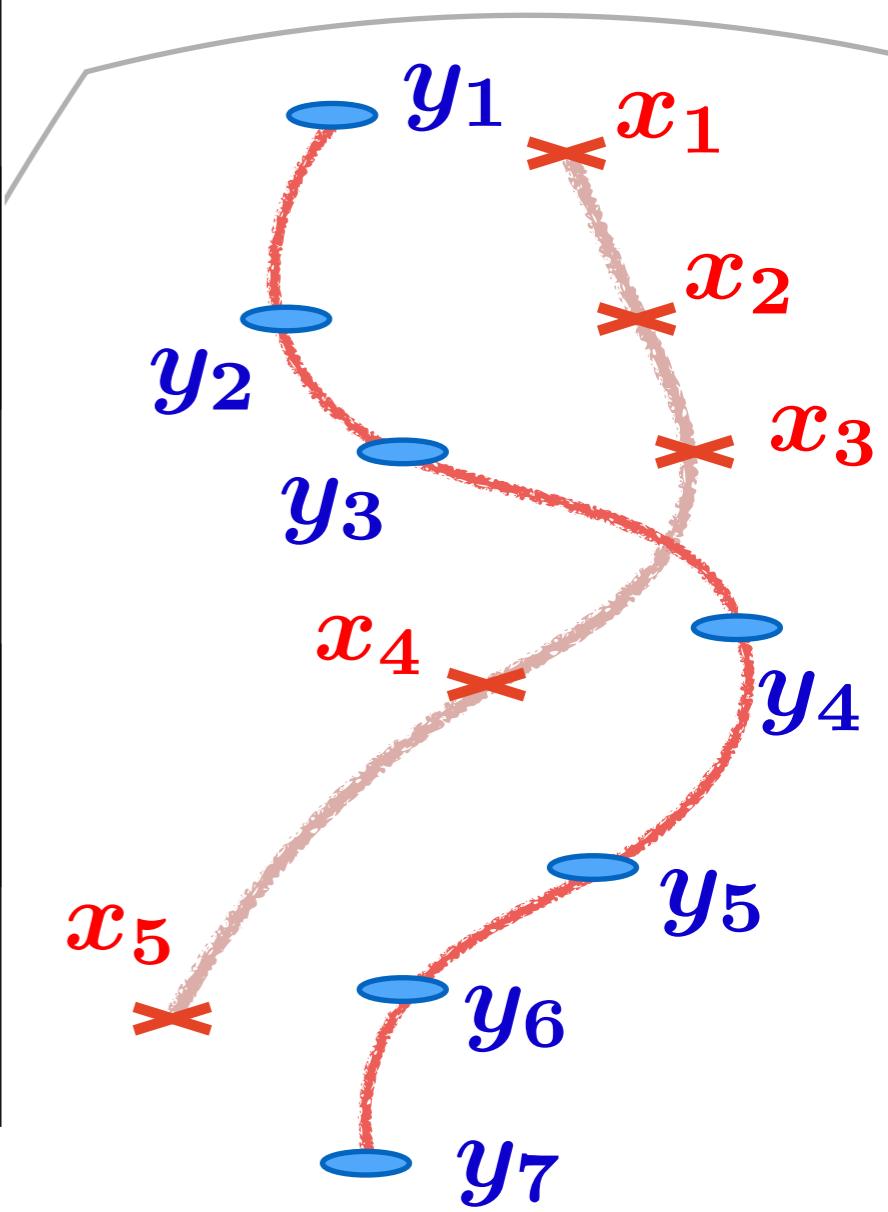
# Dynamic Time Warping [Sakoe&Chiba'78]

A discrepancy function between  
two **time series of observations**  
supported **on a metric space**.



# Pairwise Distance Matrix

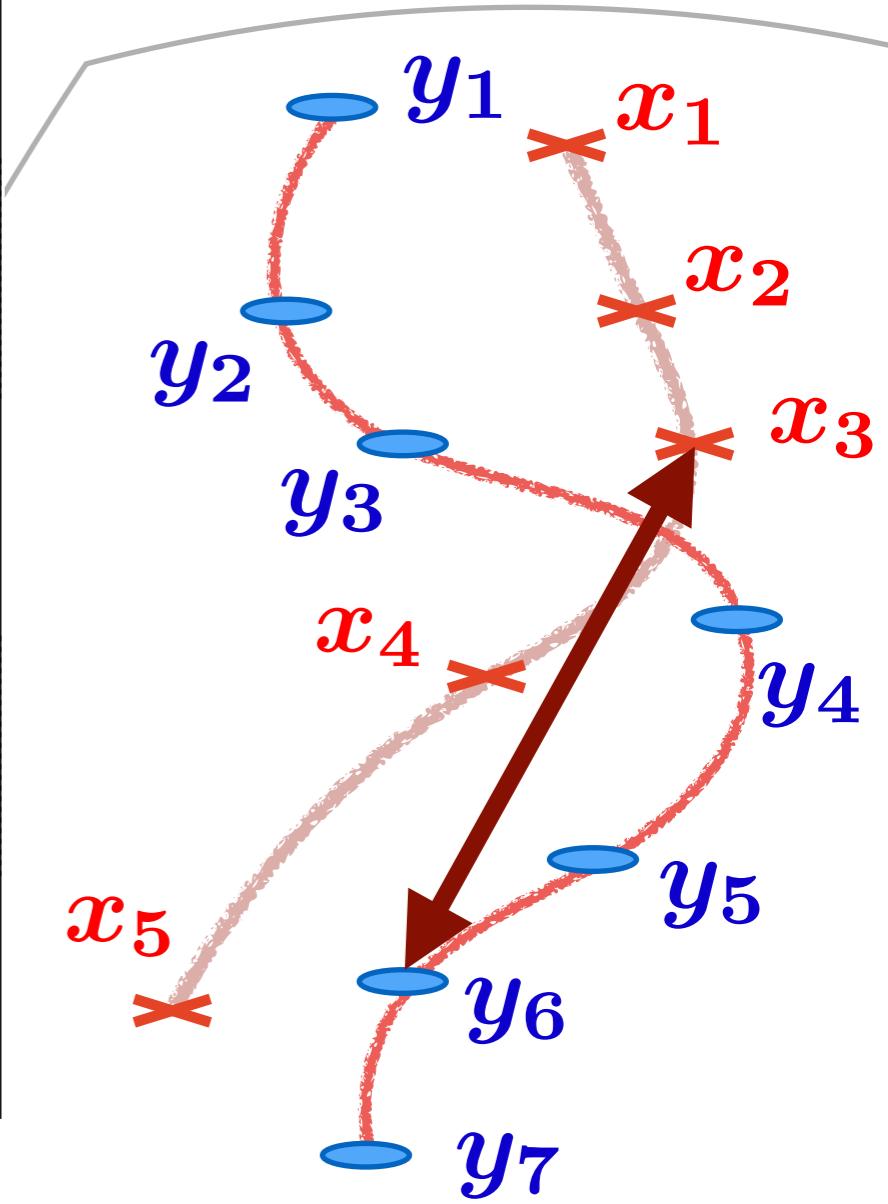
	$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$	$y_7$
$x_1$							
$x_2$							
$x_3$							
$x_4$							
$x_5$							



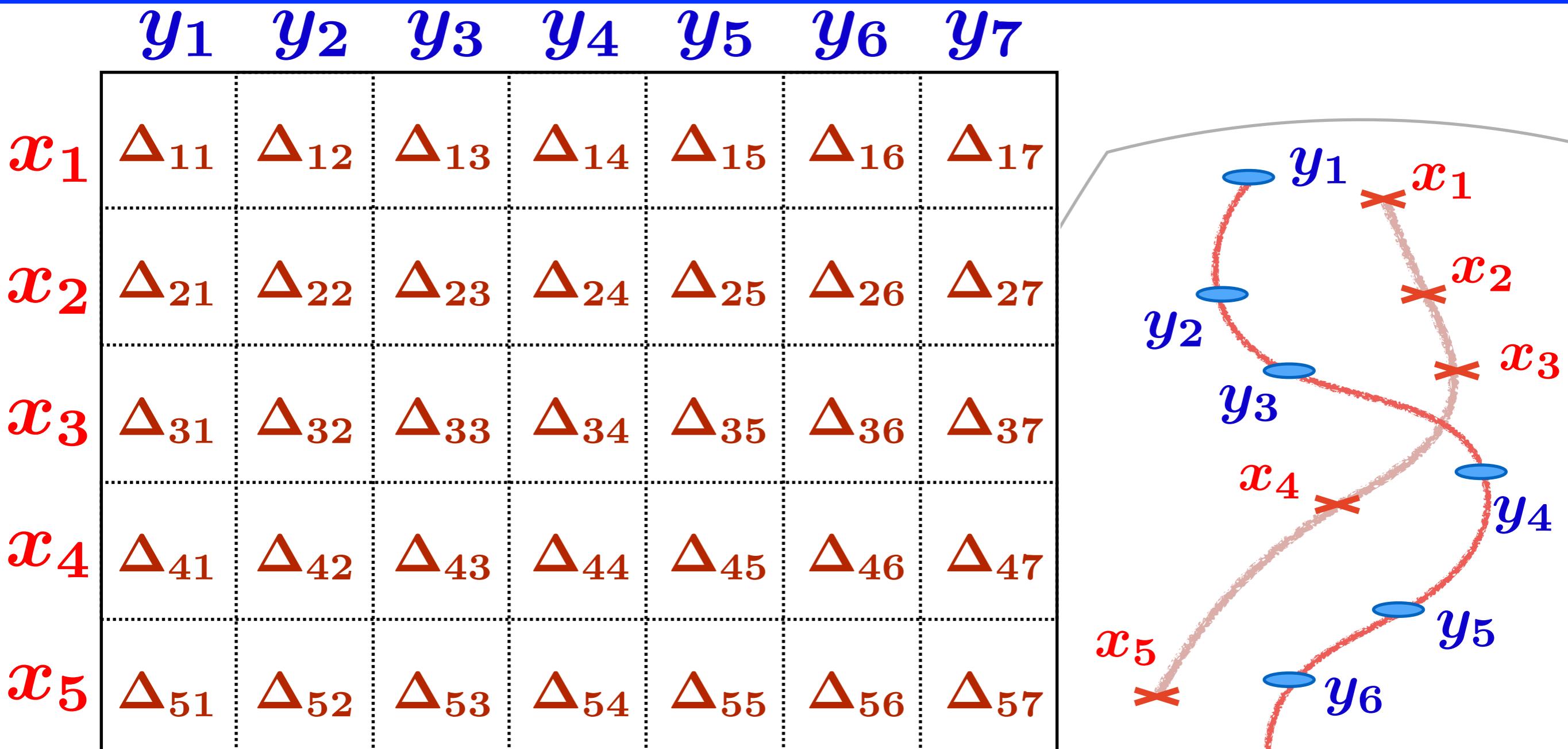
# Pairwise Distance Matrix

$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$	$y_7$
$x_1$						
$x_2$						
$x_3$						
$x_4$						
$x_5$						

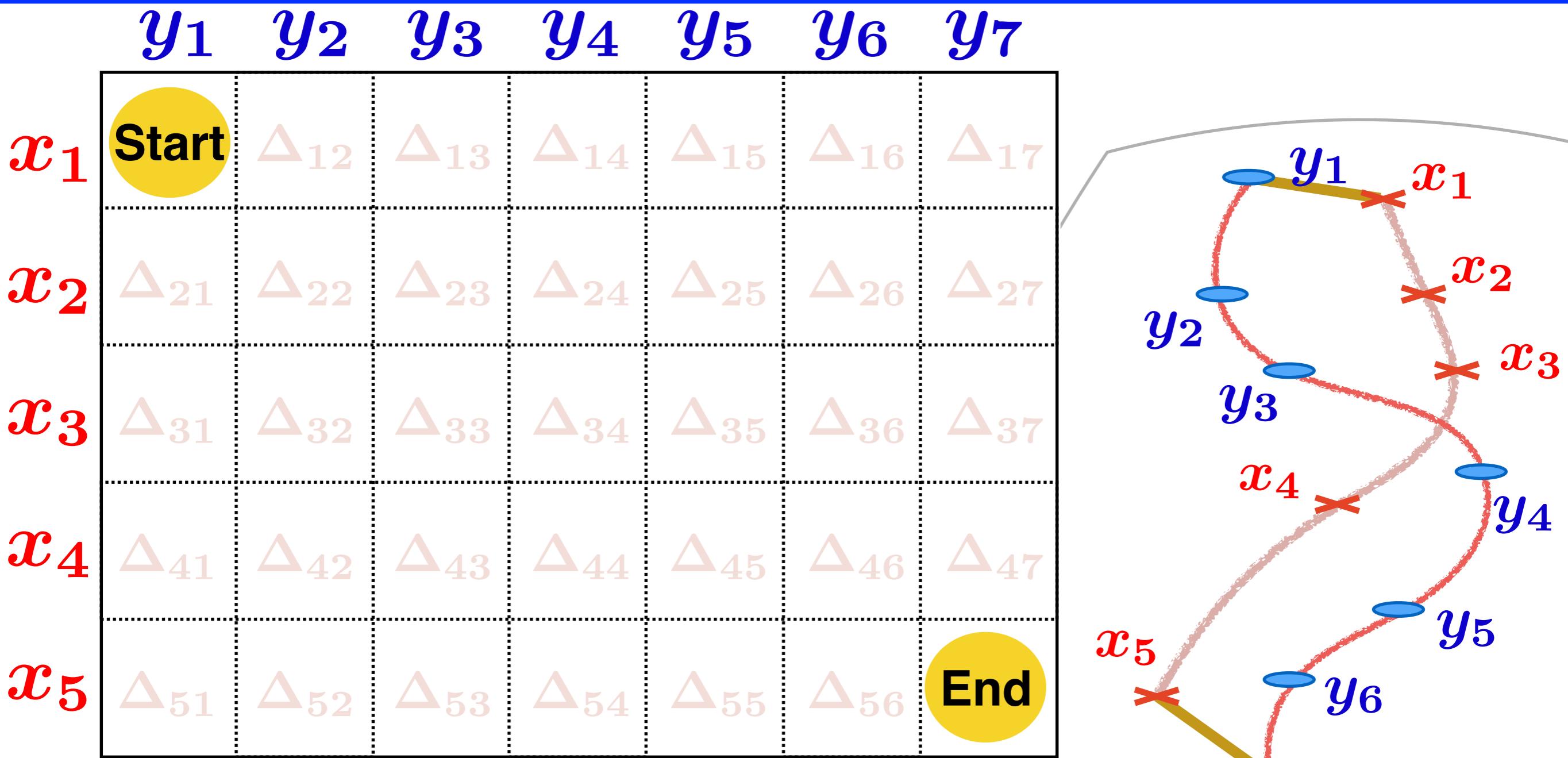
$$\Delta_{ij} = \delta(x_i, y_j)$$



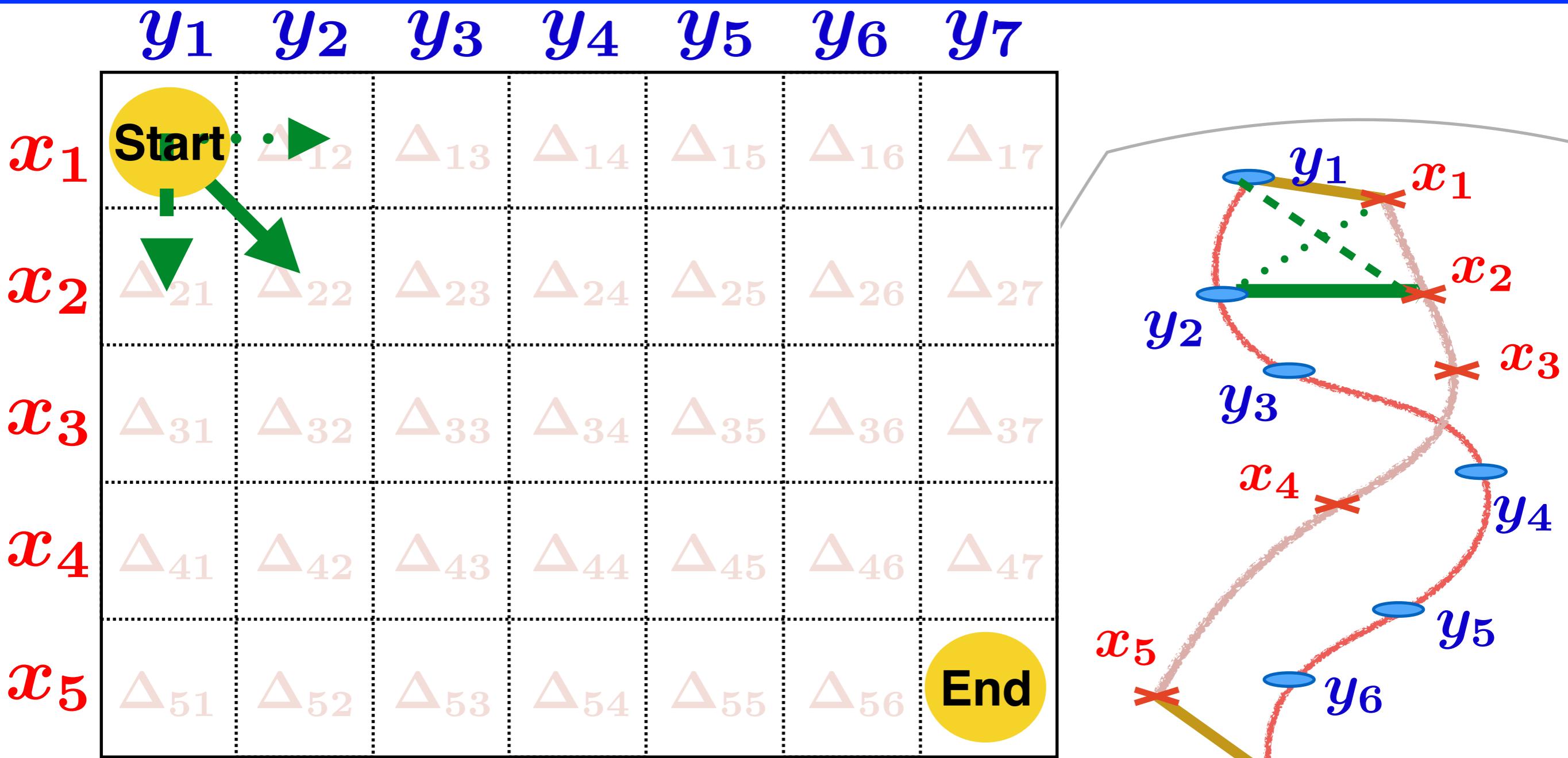
# Pairwise Distance Matrix



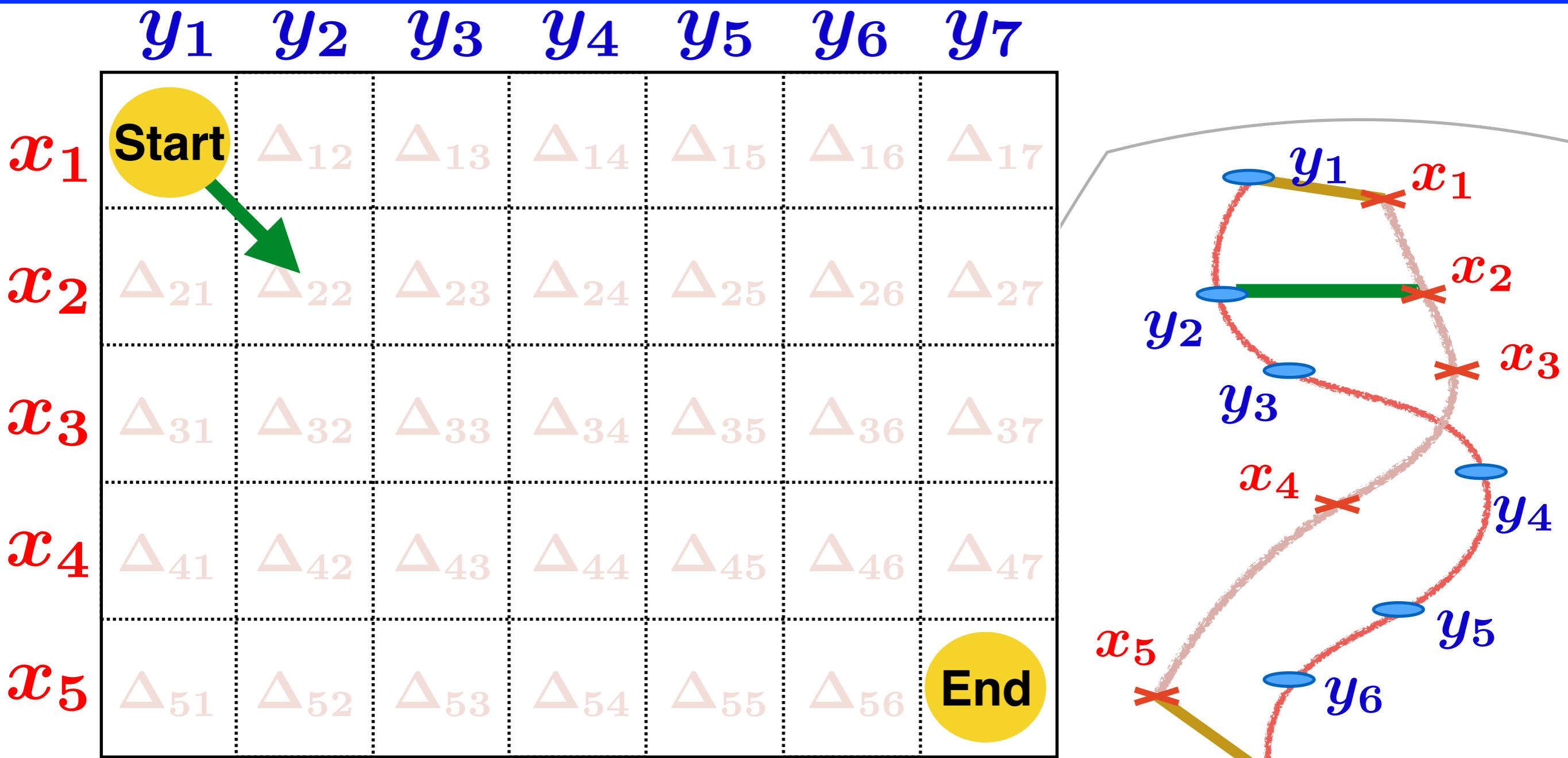
# Alignment Path



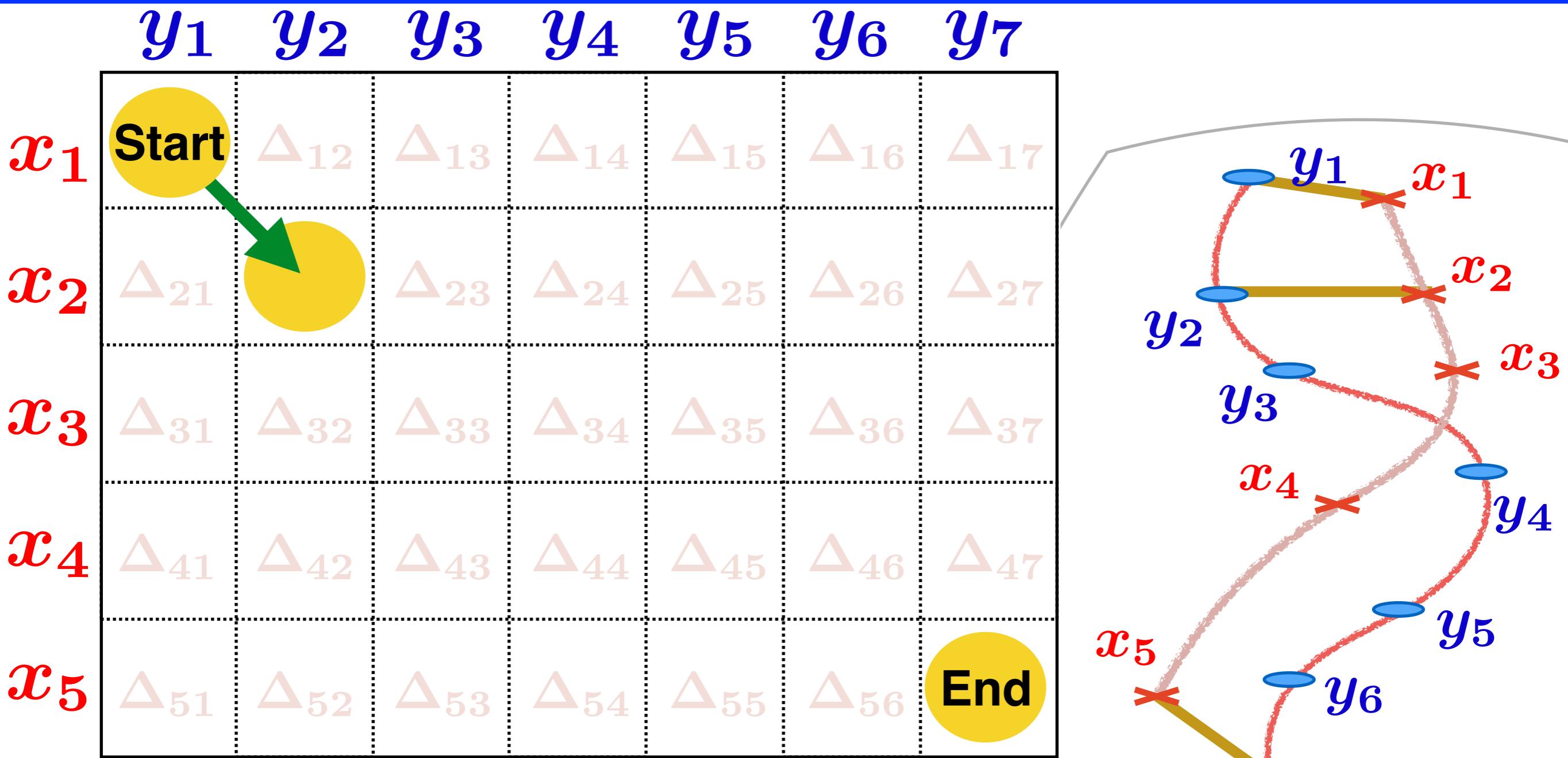
# Alignment Path



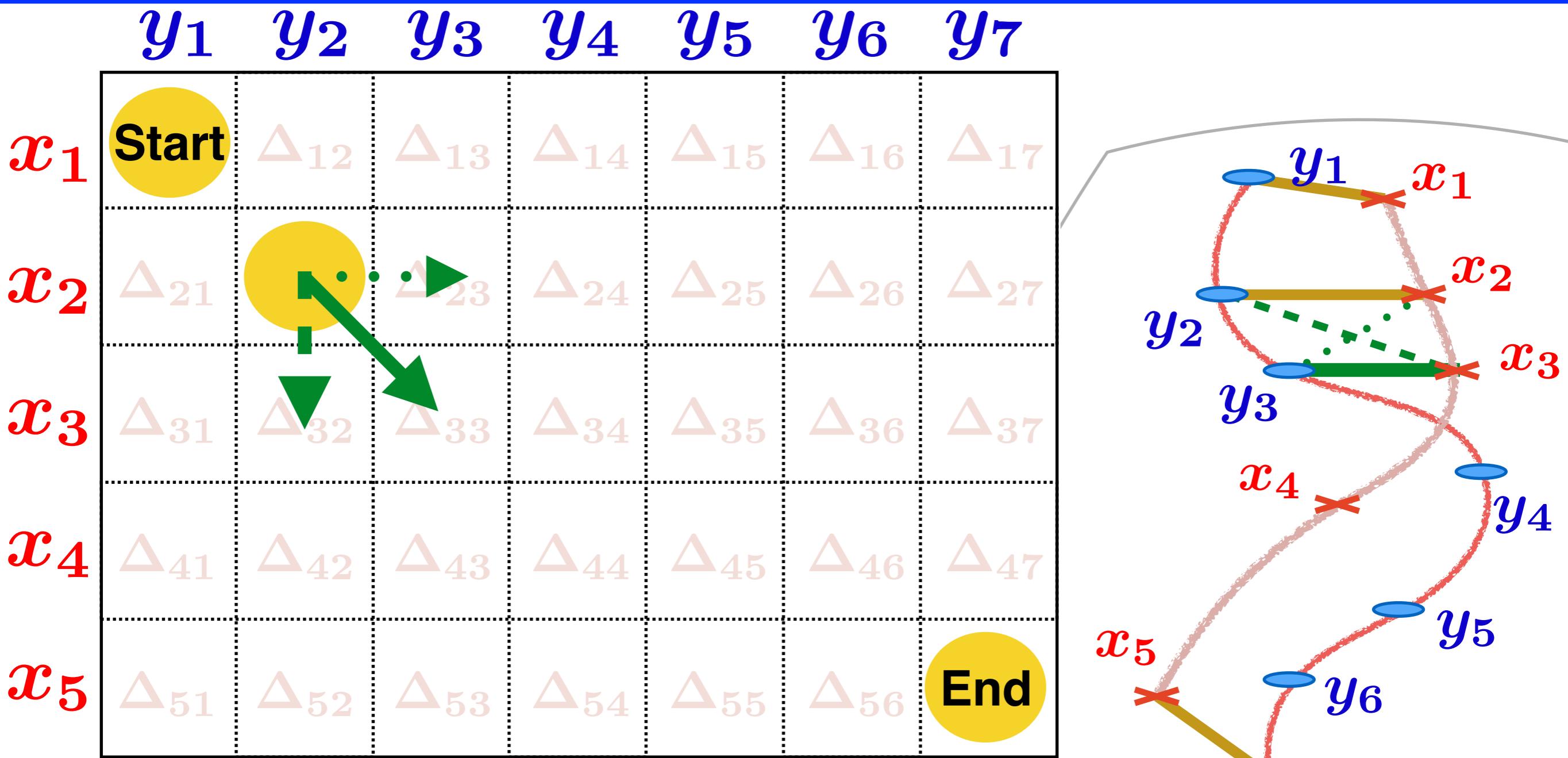
# Alignment Path



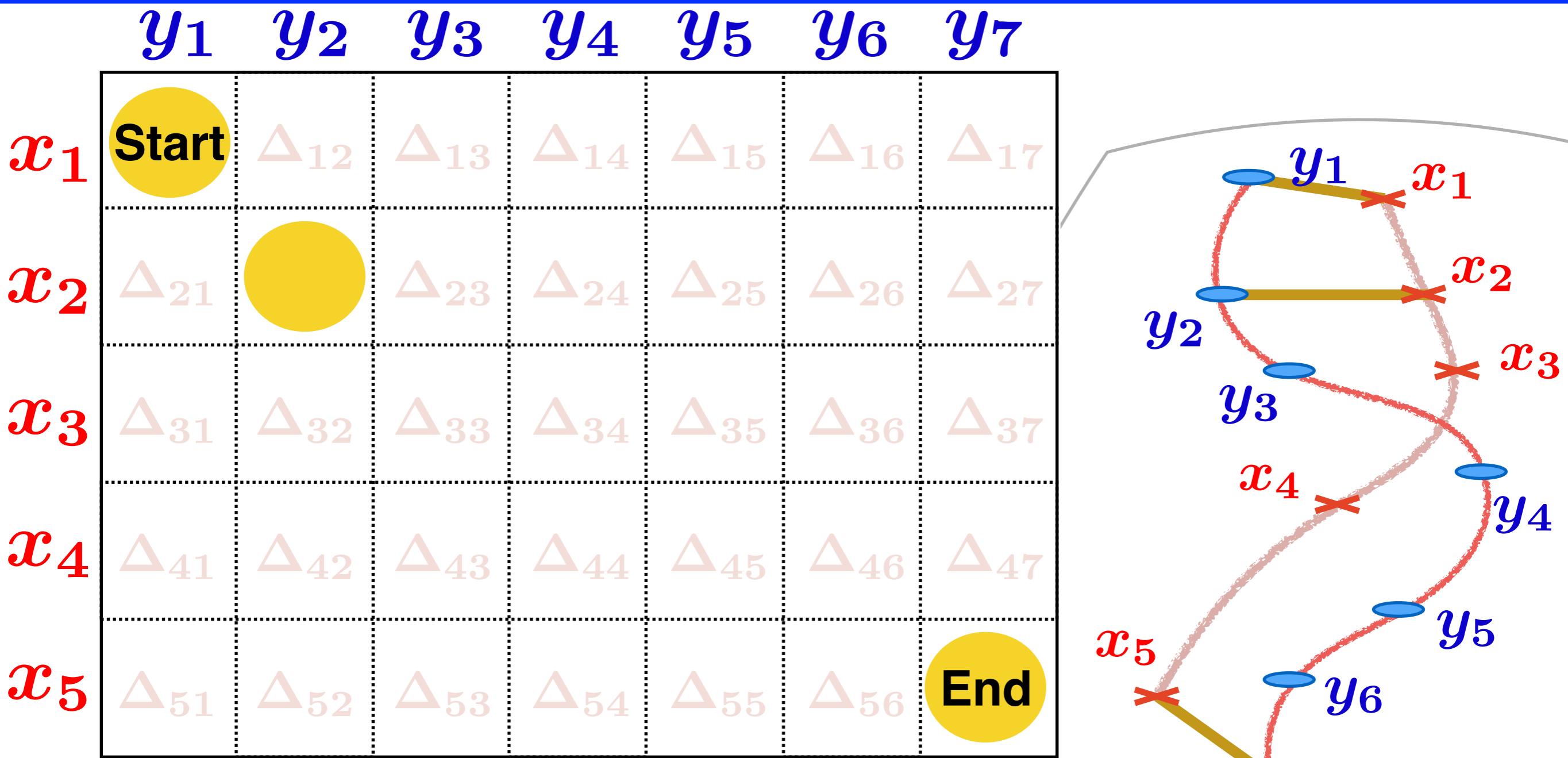
# Alignment Path



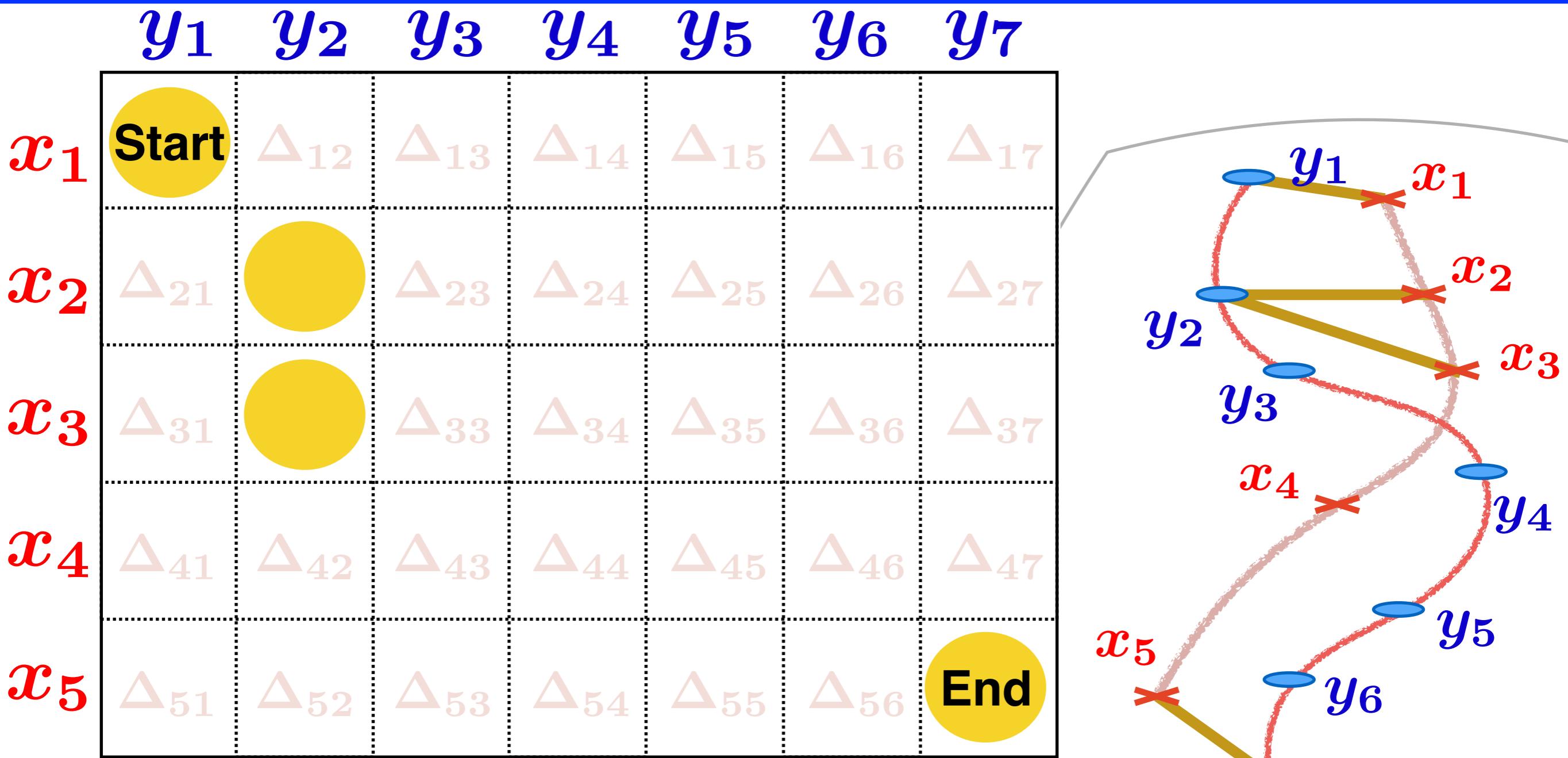
# Alignment Path



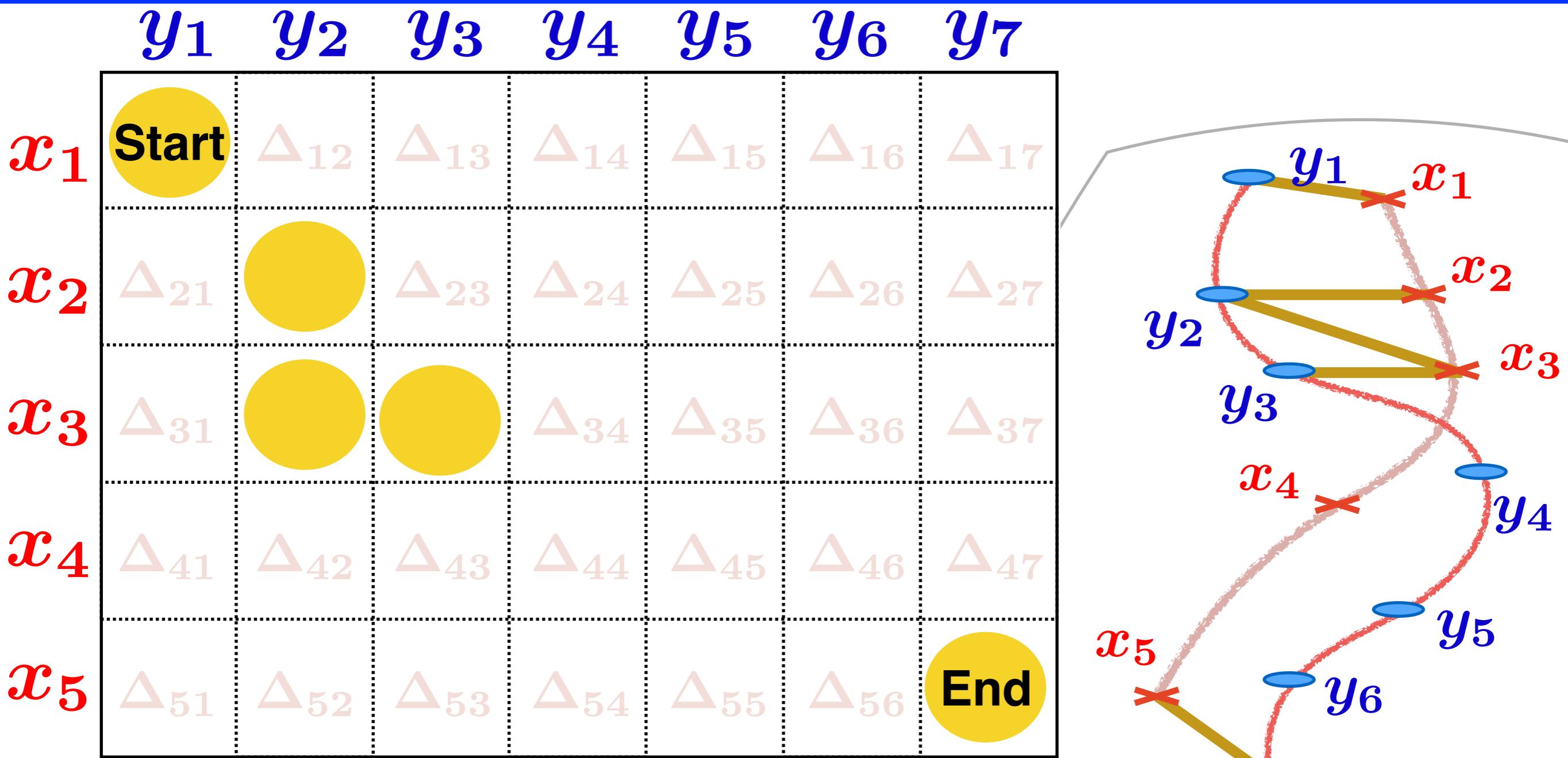
# Alignment Path



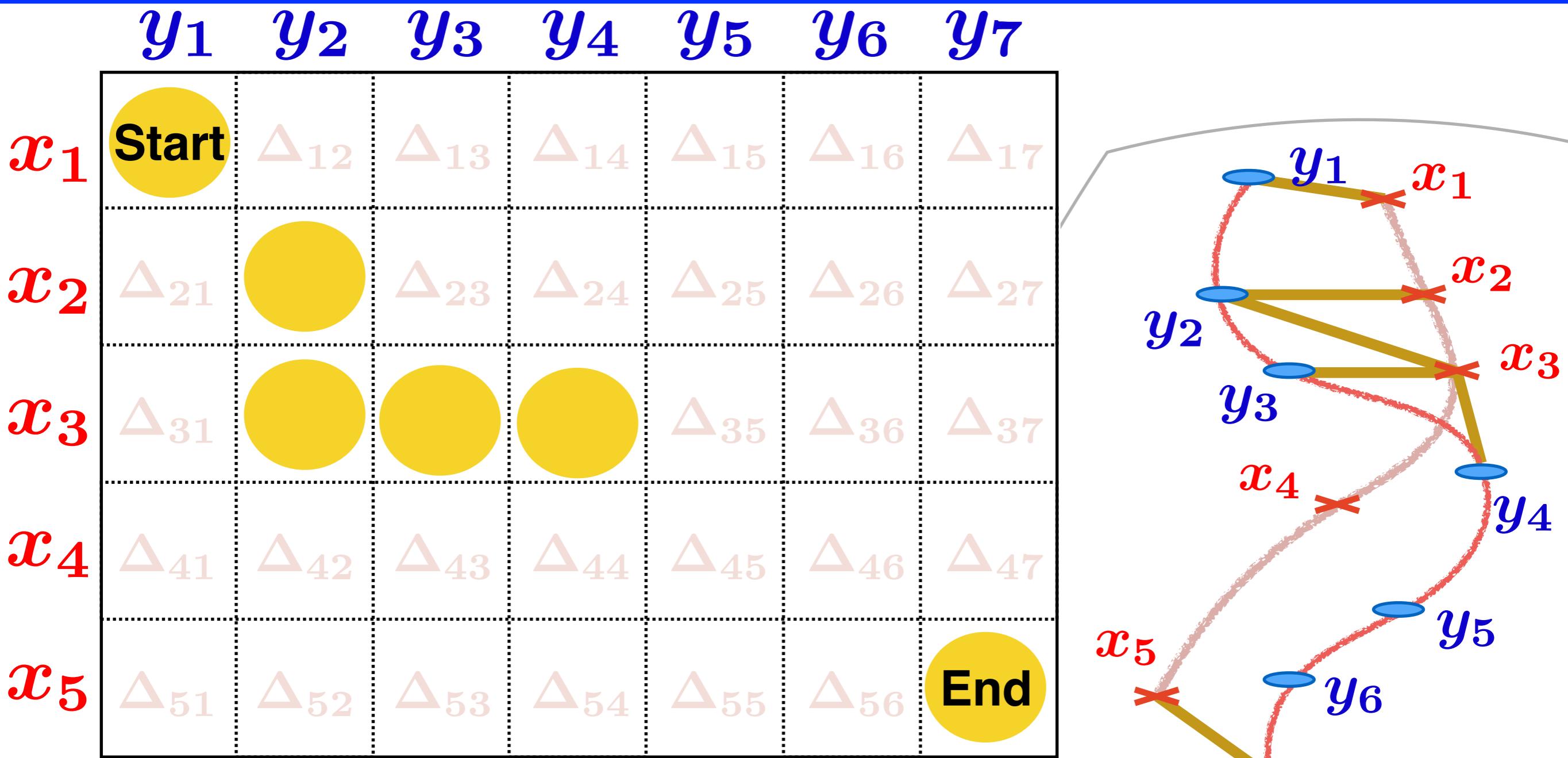
# Alignment Path



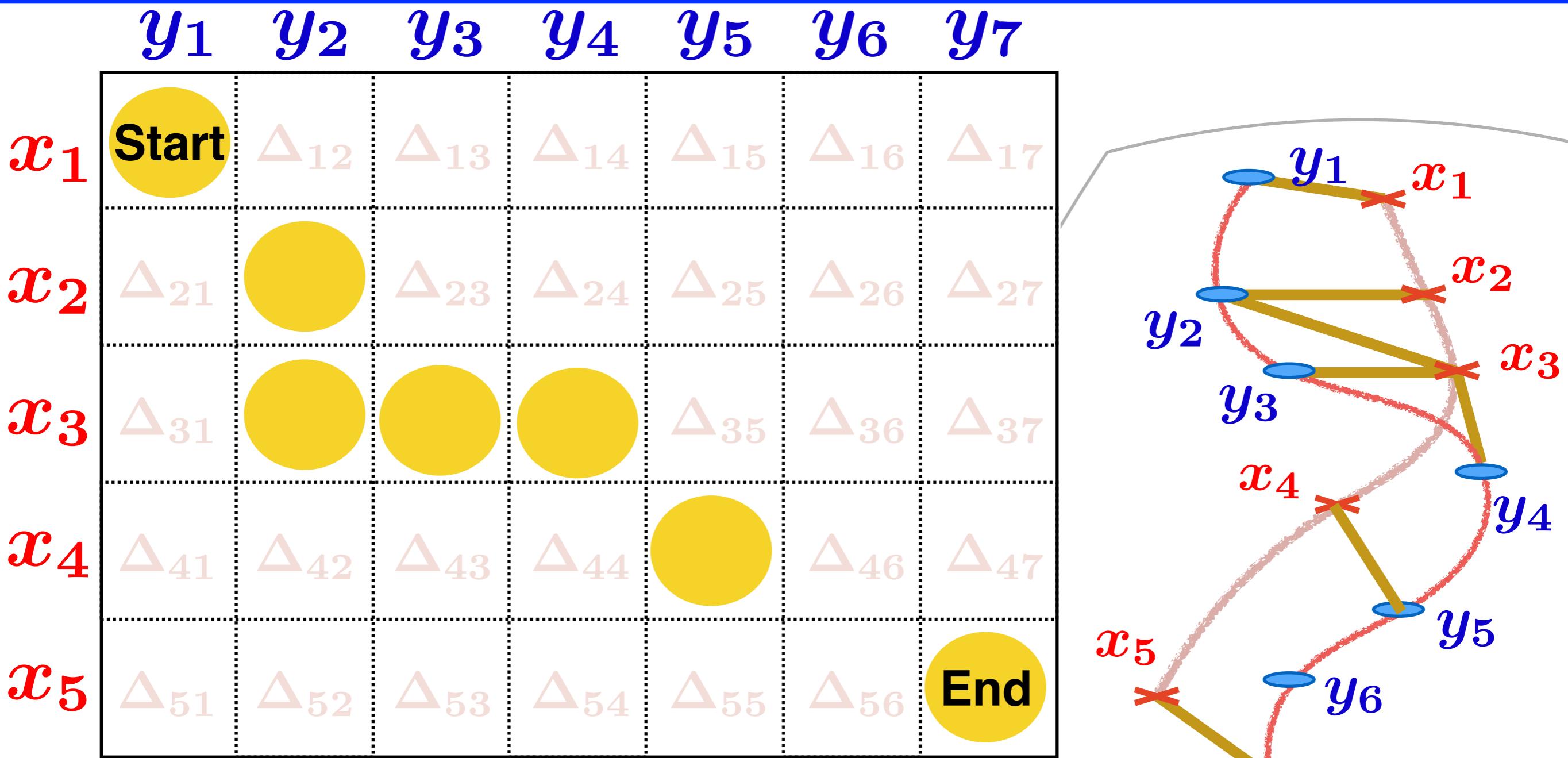
# Alignment Path



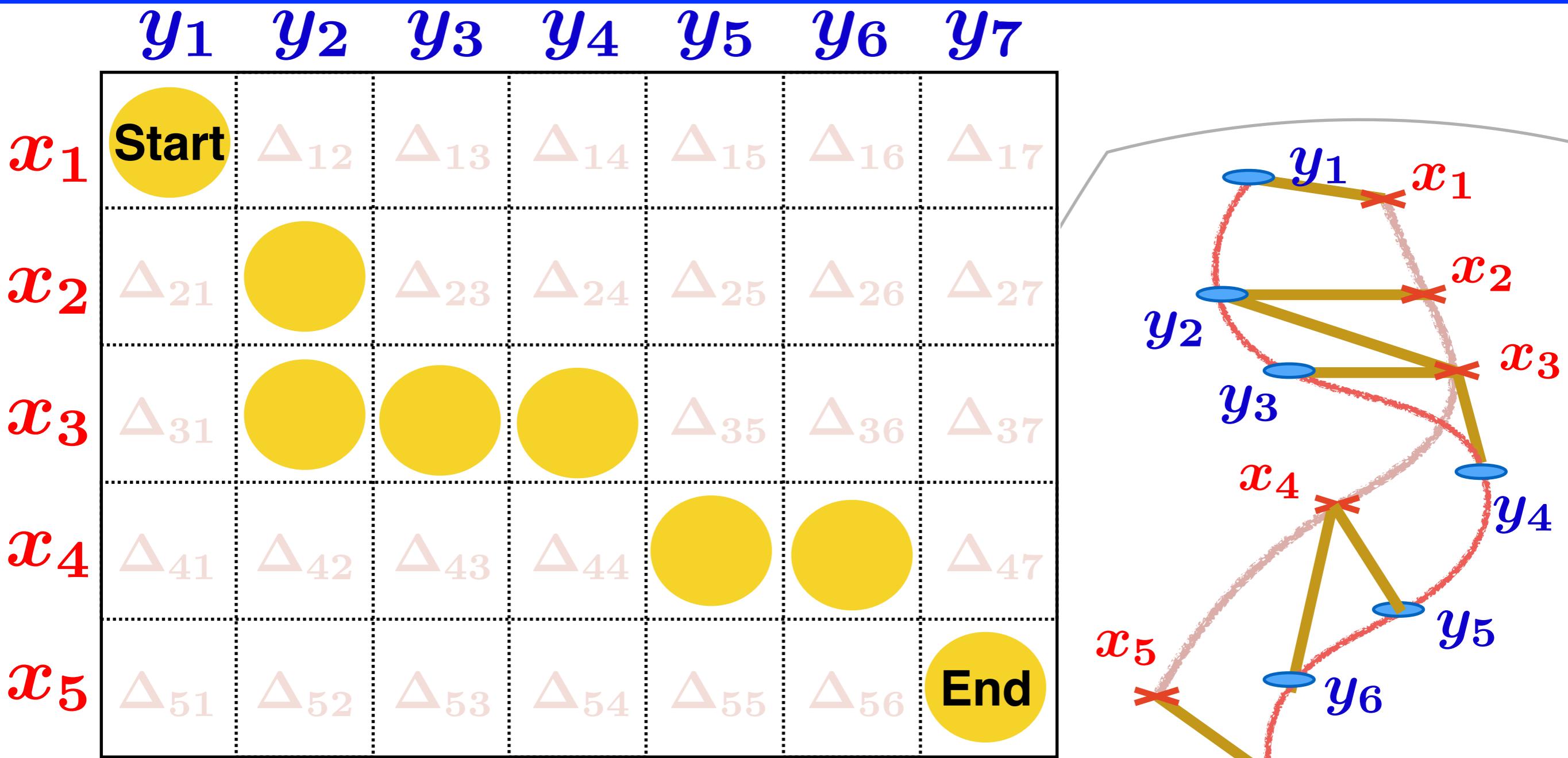
# Alignment Path



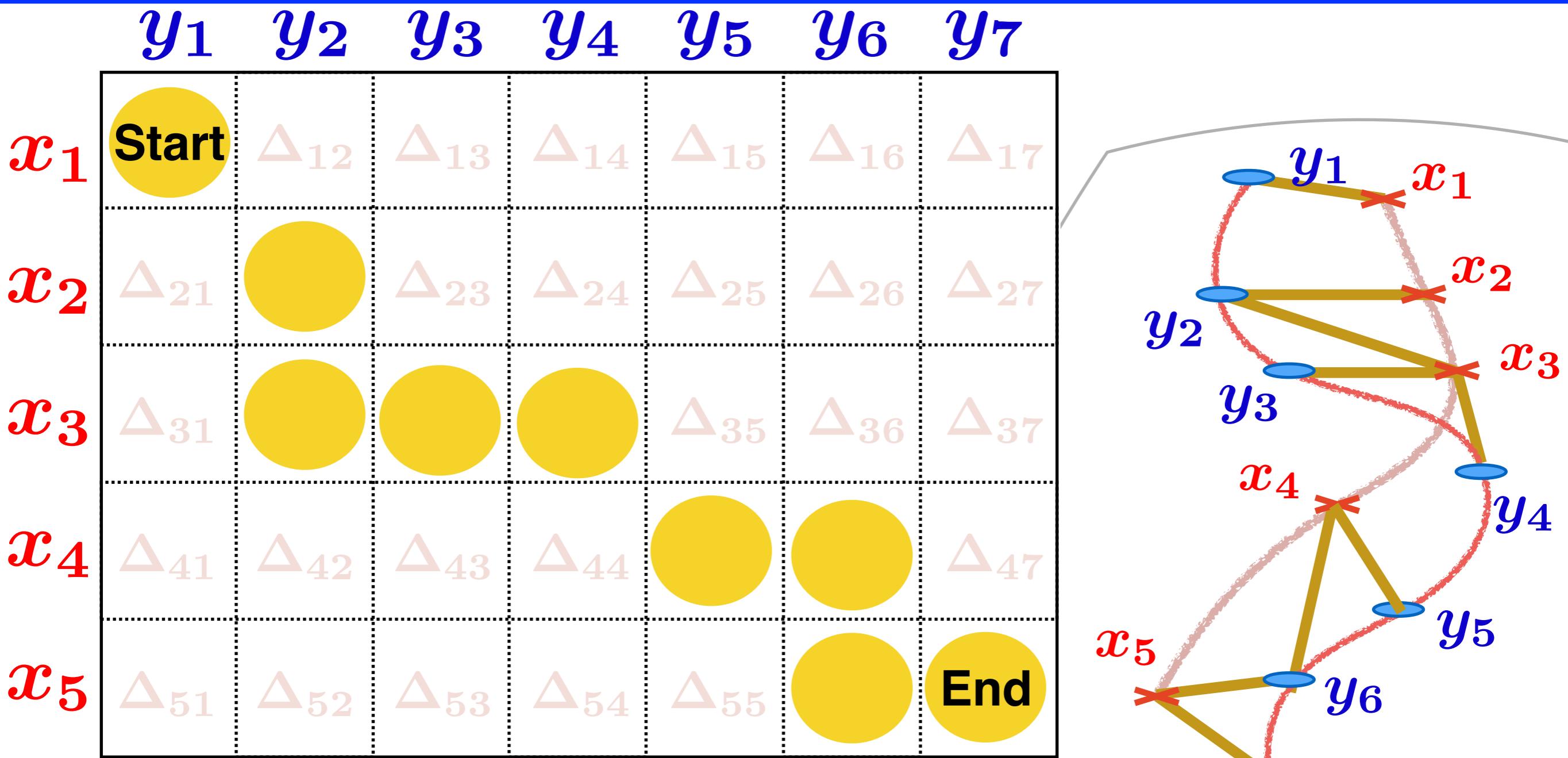
# Alignment Path



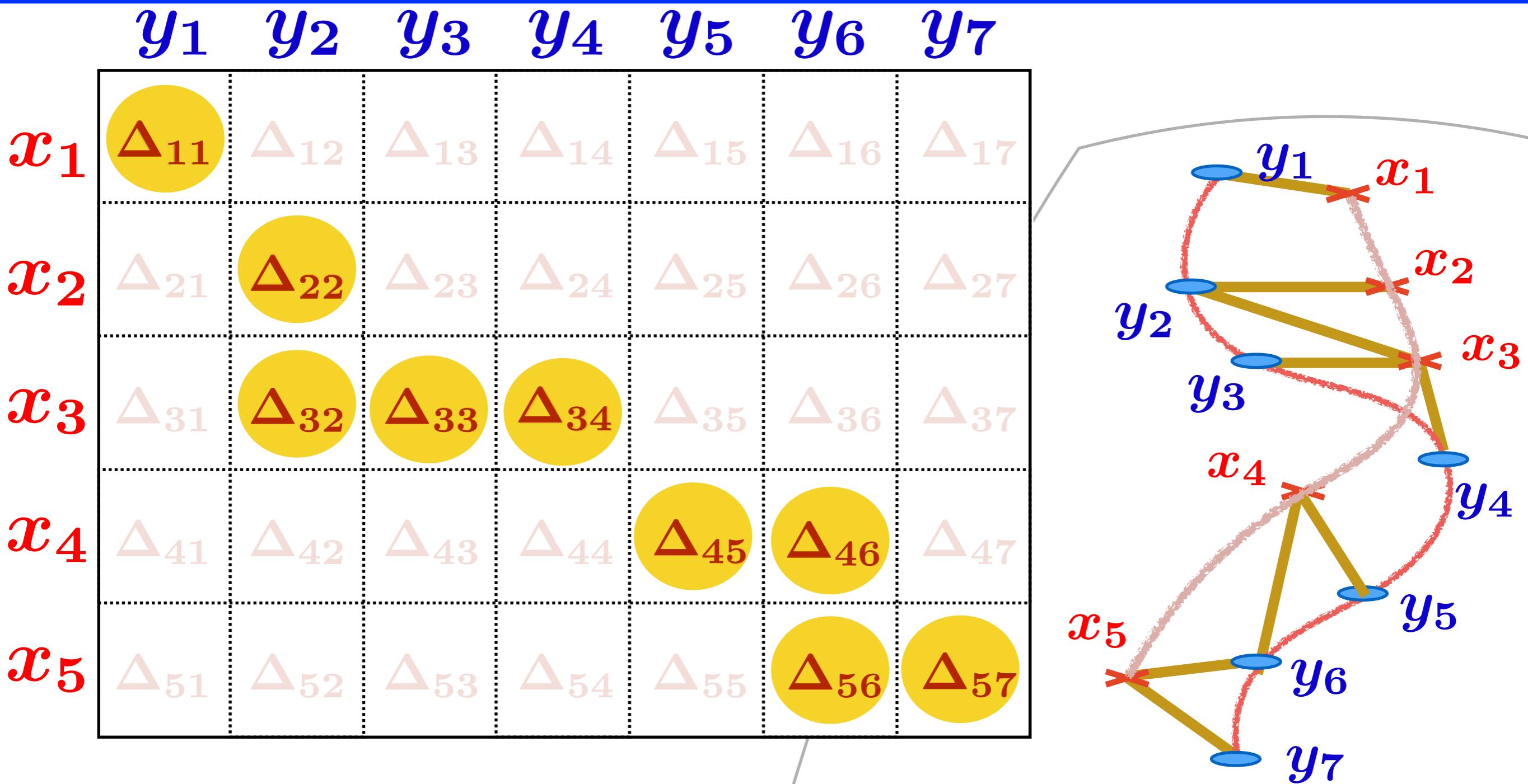
# Alignment Path



# Alignment Path

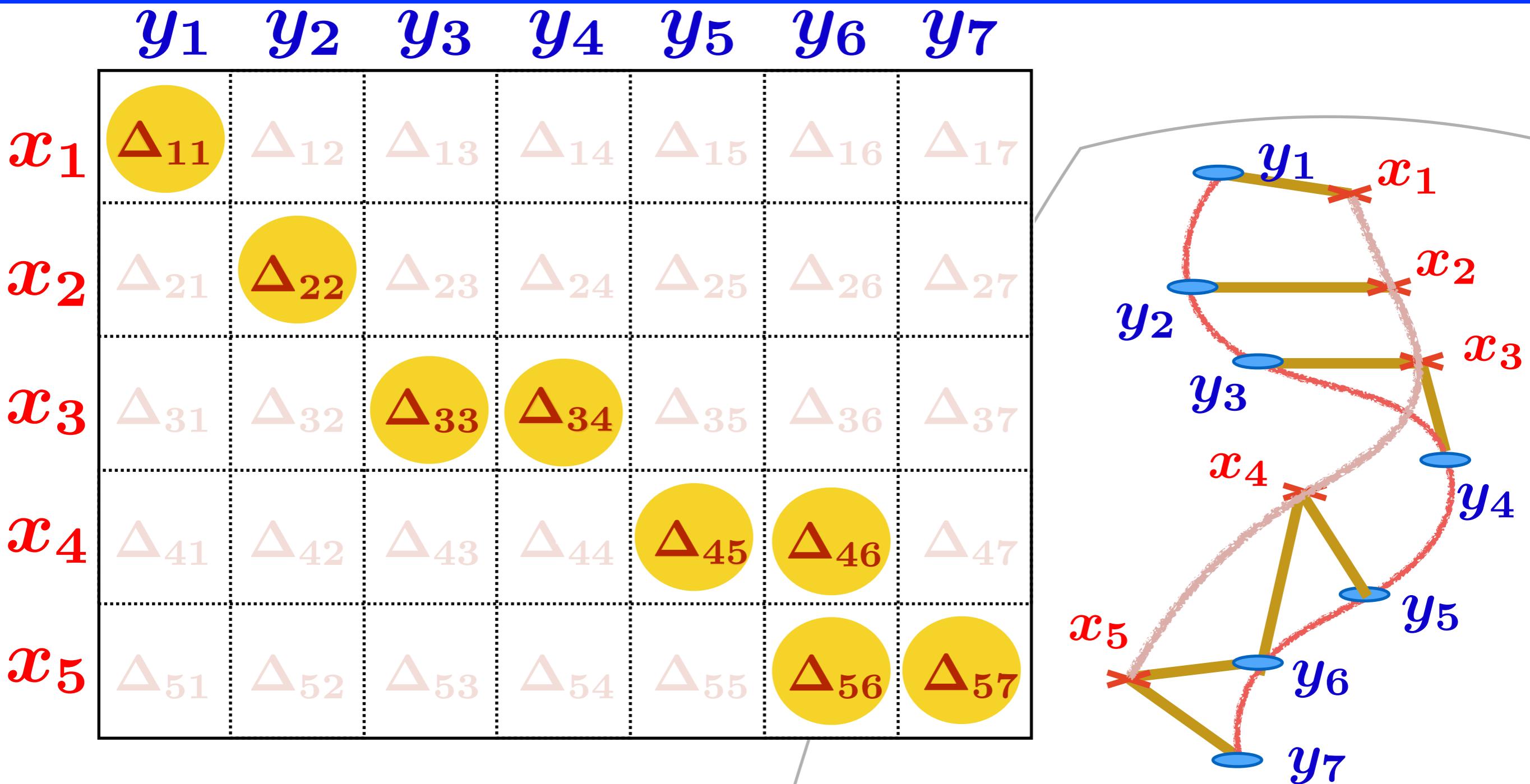


# Path Cost



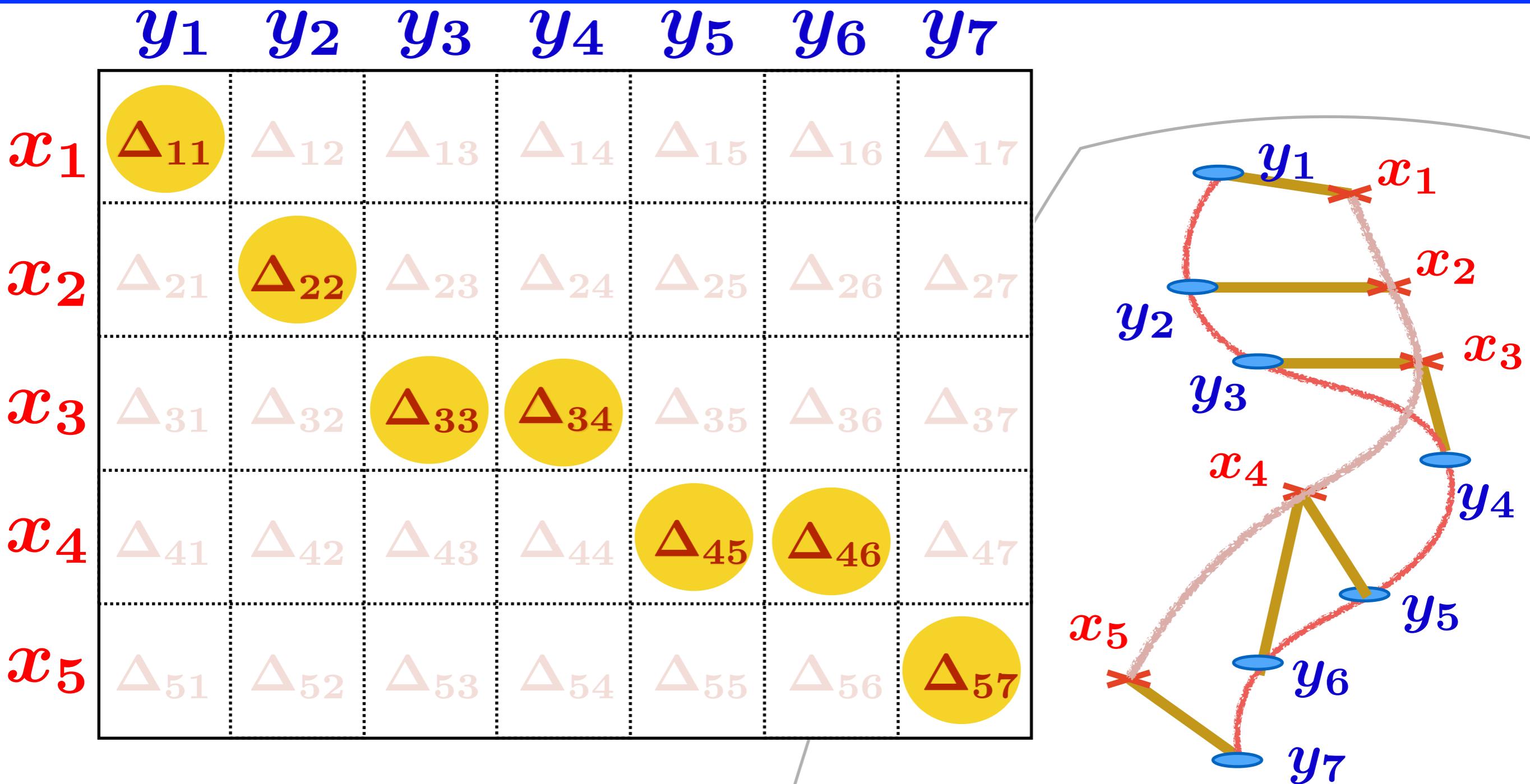
$$\text{Cost} = \Delta_{11} + \Delta_{22} + \Delta_{32} + \Delta_{33} + \Delta_{34} + \Delta_{46} + \Delta_{56} + \Delta_{57}$$

# Path Cost



$$\text{Cost} = \Delta_{11} + \Delta_{22} + \Delta_{33} + \Delta_{34} + \Delta_{46} + \Delta_{56} + \Delta_{57}$$

# Path Cost



$$\text{Cost} = \Delta_{11} + \Delta_{22} + \Delta_{33} + \Delta_{34} + \Delta_{46} + \Delta_{57}$$

# Path Cost

	$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$	$y_7$
$x_1$	$\Delta_{11}$	$\Delta_{12}$	$\Delta_{13}$	$\Delta_{14}$	$\Delta_{15}$	$\Delta_{16}$	$\Delta_{17}$
$x_2$	$\Delta_{21}$	$\Delta_{22}$	$\Delta_{23}$	$\Delta_{24}$	$\Delta_{25}$	$\Delta_{26}$	$\Delta_{27}$
$x_3$	$\Delta_{31}$	$\Delta_{32}$	$\Delta_{33}$	$\Delta_{34}$	$\Delta_{35}$	$\Delta_{36}$	$\Delta_{37}$
$x_4$	$\Delta_{41}$	$\Delta_{42}$	$\Delta_{43}$	$\Delta_{44}$	$\Delta_{45}$	$\Delta_{46}$	$\Delta_{47}$
$x_5$	$\Delta_{51}$	$\Delta_{52}$	$\Delta_{53}$	$\Delta_{54}$	$\Delta_{55}$	$\Delta_{56}$	$\Delta_{57}$

# Path Cost

	$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$	$y_7$
$x_1$	$\Delta_{11}$	$\Delta_{12}$	$\Delta_{13}$	$\Delta_{14}$	$\Delta_{15}$	$\Delta_{16}$	$\Delta_{17}$
$x_2$	$\Delta_{21}$	$\Delta_{22}$	$\Delta_{23}$	$\Delta_{24}$	$\Delta_{25}$	$\Delta_{26}$	$\Delta_{27}$
$x_3$	$\Delta_{31}$	$\Delta_{32}$	$\Delta_{33}$	$\Delta_{34}$	$\Delta_{35}$	$\Delta_{36}$	$\Delta_{37}$
$x_4$	$\Delta_{41}$	$\Delta_{42}$	$\Delta_{43}$	$\Delta_{44}$	$\Delta_{45}$	$\Delta_{46}$	$\Delta_{47}$
$x_5$	$\Delta_{51}$	$\Delta_{52}$	$\Delta_{53}$	$\Delta_{54}$	$\Delta_{55}$	$\Delta_{56}$	$\Delta_{57}$

$= A$

# Path Cost

	$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$	$y_7$
$x_1$	$\Delta_{11}$	$\Delta_{12}$	$\Delta_{13}$	$\Delta_{14}$	$\Delta_{15}$	$\Delta_{16}$	$\Delta_{17}$
$x_2$	$\Delta_{21}$	$\Delta_{22}$	$\Delta_{23}$	$\Delta_{24}$	$\Delta_{25}$	$\Delta_{26}$	$\Delta_{27}$
$x_3$	$\Delta_{31}$	$\Delta_{32}$	$\Delta_{33}$	$\Delta_{34}$	$\Delta_{35}$	$\Delta_{36}$	$\Delta_{37}$
$x_4$	$\Delta_{41}$	$\Delta_{42}$	$\Delta_{43}$	$\Delta_{44}$	$\Delta_{45}$	$\Delta_{46}$	$\Delta_{47}$
$x_5$	$\Delta_{51}$	$\Delta_{52}$	$\Delta_{53}$	$\Delta_{54}$	$\Delta_{55}$	$\Delta_{56}$	$\Delta_{57}$

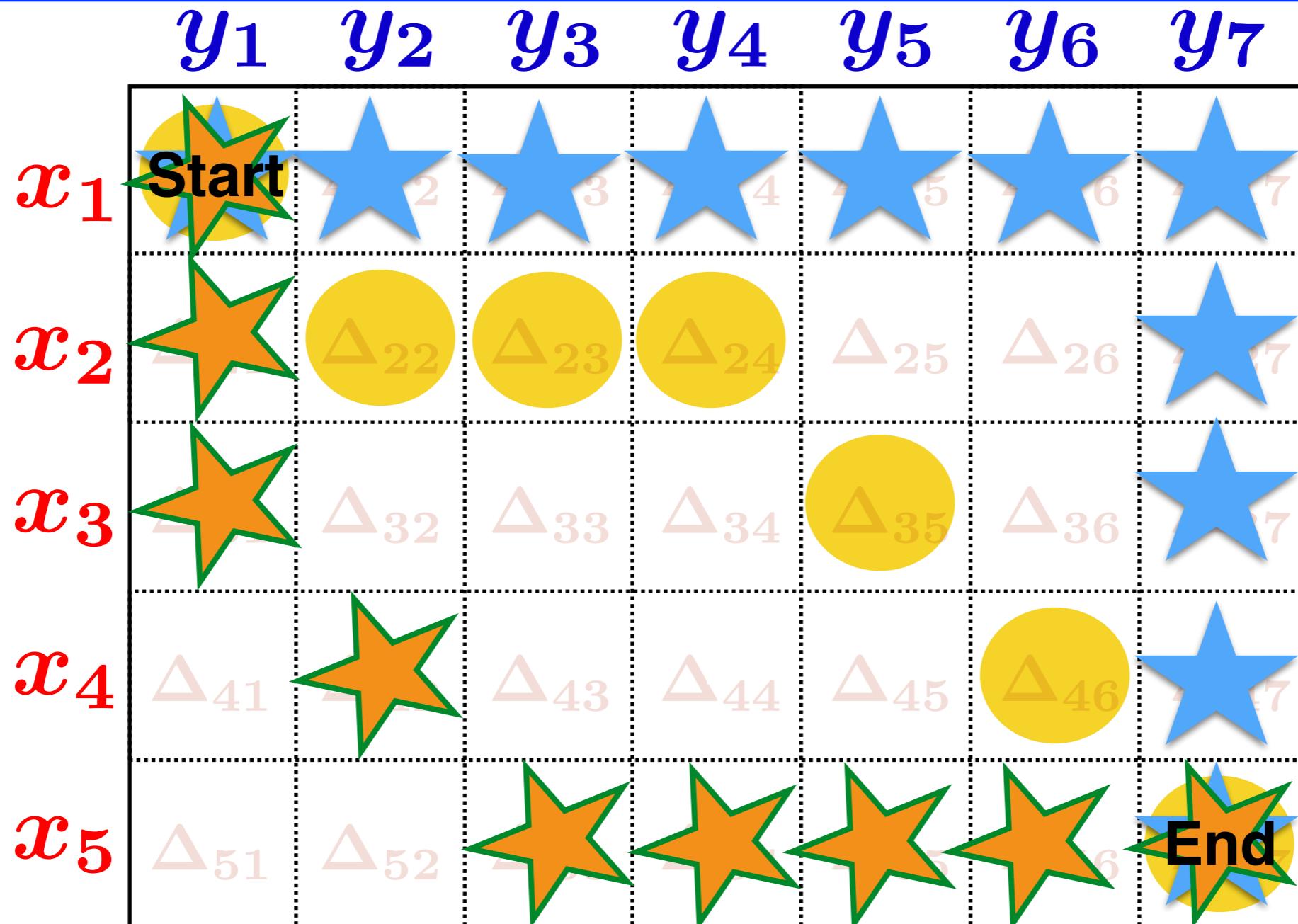
$= A$

Cost =  $\langle A, \Delta \rangle$ ,  $A \in \{0, 1\}^{n \times m}$

# Minimum Cost Alignment Matrix?

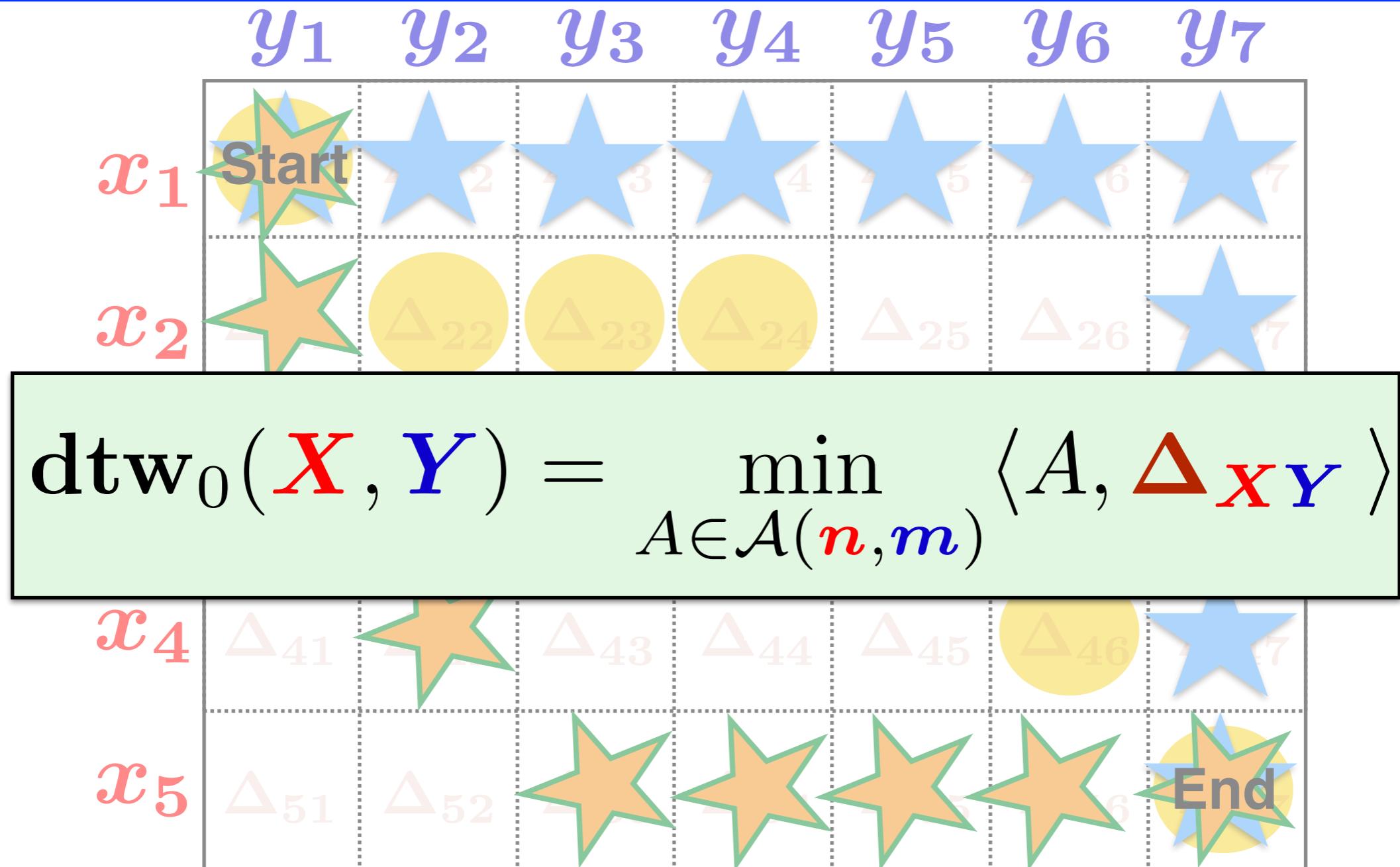
	$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$	$y_7$
$x_1$	<b>Start</b>	$\Delta_{12}$	$\Delta_{13}$	$\Delta_{14}$	$\Delta_{15}$	$\Delta_{16}$	$\Delta_{17}$
$x_2$	$\Delta_{21}$	$\Delta_{22}$	$\Delta_{23}$	$\Delta_{24}$	$\Delta_{25}$	$\Delta_{26}$	$\Delta_{27}$
$x_3$	$\Delta_{31}$	$\Delta_{32}$	$\Delta_{33}$	$\Delta_{34}$	$\Delta_{35}$	$\Delta_{36}$	$\Delta_{37}$
$x_4$	$\Delta_{41}$	$\Delta_{42}$	$\Delta_{43}$	$\Delta_{44}$	$\Delta_{45}$	$\Delta_{46}$	$\Delta_{47}$
$x_5$	$\Delta_{51}$	$\Delta_{52}$	$\Delta_{53}$	$\Delta_{54}$	$\Delta_{55}$	$\Delta_{56}$	<b>End</b>

# Minimum Cost Alignment Matrix?



Set of all valid path matrices:  $\mathcal{A}(n, m) \subset \{0, 1\}^{n \times m}$

# Dynamic Time Warping [Sakoe&Chiba'78]



Set of all valid path matrices:  $\mathcal{A}(n, m) \subset \{0, 1\}^{n \times m}$

# Number of valid paths

Size of  $\mathcal{A}(n, m)$  is exponential in  $n, m$ .

$$\#\mathcal{A}(n, m) = \text{Delannoy}(n - 1, m - 1)$$

n=3, m=3	321
n=5, m=5	1683
n=10, m=10	8097453
...	

Set of all valid path matrices:  $\mathcal{A}(n, m) \subset \{0, 1\}^{n \times m}$

# Best Path: Bellman Recursion

	$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$	$y_7$
$x_1$	$\Delta_{11}$	$\Delta_{12}$	$\Delta_{13}$	$\Delta_{14}$	$\Delta_{15}$	$\Delta_{16}$	$\Delta_{17}$
$x_2$	$\Delta_{21}$	$\Delta_{22}$	$\Delta_{23}$	$\Delta_{24}$	$\Delta_{25}$	$\Delta_{26}$	$\Delta_{27}$
$x_3$	$\Delta_{31}$	$\Delta_{32}$	$\Delta_{33}$	$\Delta_{34}$	$\Delta_{35}$	$\Delta_{36}$	$\Delta_{37}$
$x_4$	$\Delta_{41}$	$\Delta_{42}$	$\Delta_{43}$	$\Delta_{44}$	$\Delta_{45}$	$\Delta_{46}$	$\Delta_{47}$
$x_5$	$\Delta_{51}$	$\Delta_{52}$	$\Delta_{53}$	$\Delta_{54}$	$\Delta_{55}$	$\Delta_{56}$	$\Delta_{57}$

# Best Path: Bellman Recursion

	$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$	$y_7$
$x_1$	$\Delta_{11}$	$\Delta_{12}$	$\Delta_{13}$	$\Delta_{14}$	$\Delta_{15}$	$\Delta_{16}$	$\Delta_{17}$
$x_2$	$\Delta_{21}$	$\Delta_{22}$	$\Delta_{23}$	$\Delta_{24}$	$\Delta_{25}$	$\Delta_{26}$	$\Delta_{27}$
$x_3$	$\Delta_{31}$	$\Delta_{32}$	$\Delta_{33}$	$\Delta_{34}$	$r_{3,5}^*$	$\Delta_{36}$	$\Delta_{37}$
$x_4$	$\Delta_{41}$	$\Delta_{42}$	$\Delta_{43}$	$\Delta_{44}$	$\Delta_{45}$	$\Delta_{46}$	$\Delta_{47}$
$x_5$	$\Delta_{51}$	$\Delta_{52}$	$\Delta_{53}$	$\Delta_{54}$	$\Delta_{55}$	$\Delta_{56}$	$\Delta_{57}$

$$r_{3,5}^* = \min_{A \in \mathcal{A}(3,5)} \langle A, [\Delta_{ij}]_{i \leq 3, j \leq 5} \rangle$$

# Best Path: Bellman Recursion

	$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$	$y_7$
$x_1$	$\Delta_{11}$	$\Delta_{12}$	$\Delta_{13}$	$\Delta_{14}$	$\Delta_{15}$	$\Delta_{16}$	$\Delta_{17}$
$x_2$	$\Delta_{21}$	$\Delta_{22}$	$\Delta_{23}$	$\Delta_{24}$	$\Delta_{25}$	$\Delta_{26}$	$\Delta_{27}$
$x_3$	$\Delta_{31}$	$\Delta_{32}$	$\Delta_{33}$	$\Delta_{34}$	$r_{3,5}^*$	$\Delta_{36}$	$\Delta_{37}$
$x_4$	$\Delta_{41}$	$\Delta_{42}$	$\Delta_{43}$	$r_{4,4}^*$	$\Delta_{45}$	$\Delta_{46}$	$\Delta_{47}$
$x_5$	$\Delta_{51}$	$\Delta_{52}$	$\Delta_{53}$	$\Delta_{54}$	$\Delta_{55}$	$\Delta_{56}$	$\Delta_{57}$

# Best Path: Bellman Recursion

	$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$	$y_7$
$x_1$	$\Delta_{11}$	$\Delta_{12}$	$\Delta_{13}$	$\Delta_{14}$	$\Delta_{15}$	$\Delta_{16}$	$\Delta_{17}$
$x_2$	$\Delta_{21}$	$\Delta_{22}$	$\Delta_{23}$	$\Delta_{24}$	$\Delta_{25}$	$\Delta_{26}$	$\Delta_{27}$
$x_3$	$\Delta_{31}$	$\Delta_{32}$	$\Delta_{33}$	$r_{3,4}^*$	$r_{3,5}^*$	$\Delta_{36}$	$\Delta_{37}$
$x_4$	$\Delta_{41}$	$\Delta_{42}$	$\Delta_{43}$	$r_{4,4}^*$	$\Delta_{45}$	$\Delta_{46}$	$\Delta_{47}$
$x_5$	$\Delta_{51}$	$\Delta_{52}$	$\Delta_{53}$	$\Delta_{54}$	$\Delta_{55}$	$\Delta_{56}$	$\Delta_{57}$

# Best Path: Bellman Recursion

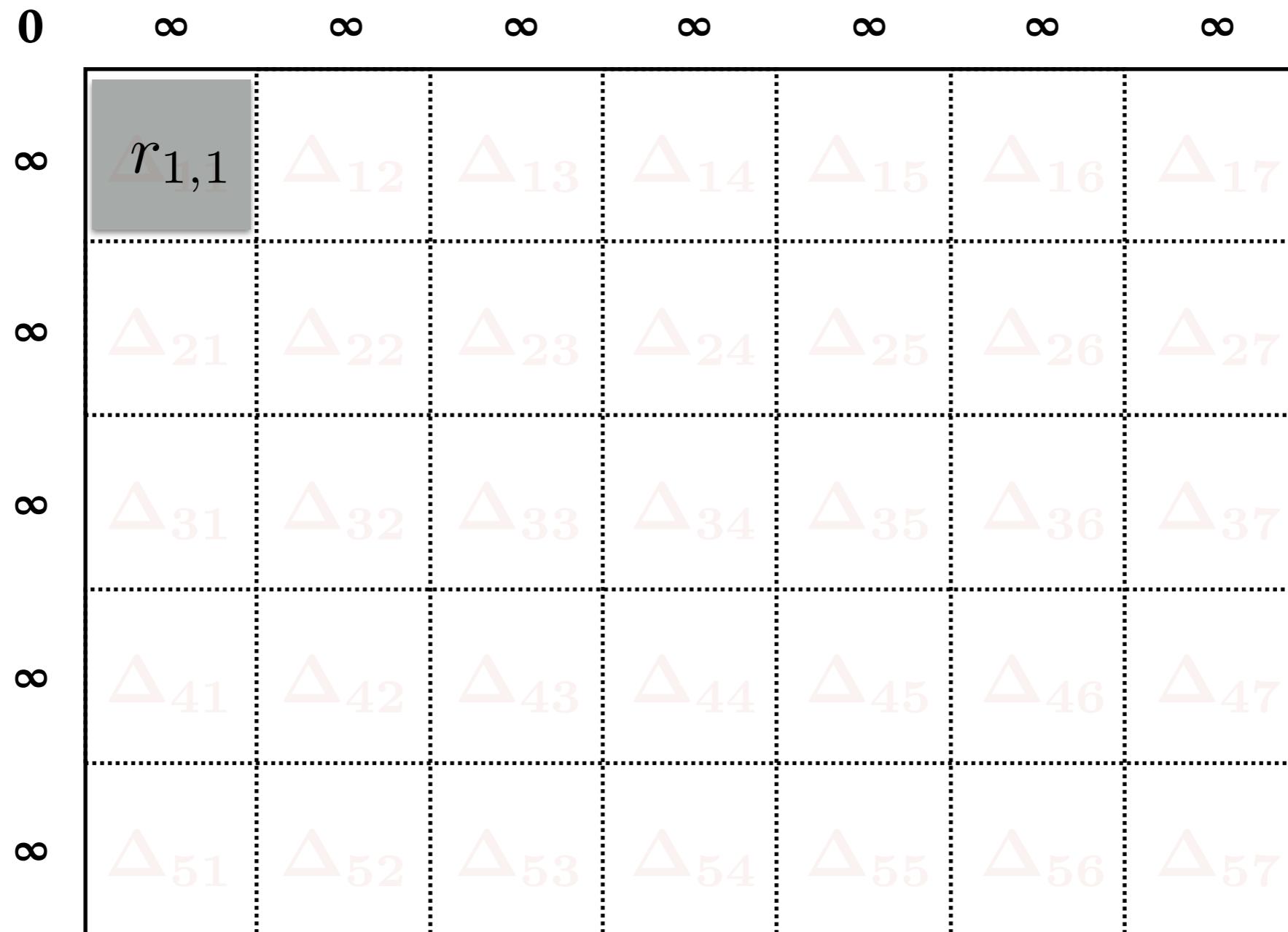
	$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$	$y_7$
$x_1$	$\Delta_{11}$	$\Delta_{12}$	$\Delta_{13}$	$\Delta_{14}$	$\Delta_{15}$	$\Delta_{16}$	$\Delta_{17}$
$x_2$	$\Delta_{21}$	$\Delta_{22}$	$\Delta_{23}$	$\Delta_{24}$	$\Delta_{25}$	$\Delta_{26}$	$\Delta_{27}$
$x_3$	$\Delta_{31}$	$\Delta_{32}$	$\Delta_{33}$	$r_{3,4}^*$	$r_{3,5}^*$	$\Delta_{36}$	$\Delta_{37}$
$x_4$	$\Delta_{41}$	$\Delta_{42}$	$\Delta_{43}$	$r_{4,4}^*$	$r_{4,5}^*$	$\Delta_{46}$	$\Delta_{47}$
$x_5$	$\Delta_{51}$	$\Delta_{52}$	$\Delta_{53}$	$\Delta_{54}$	$\Delta_{55}$	$\Delta_{56}$	$\Delta_{57}$

# Best Path: Bellman Recursion

	$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$	$y_7$
$x_1$	$\Delta_{11}$	$\Delta_{12}$	$\Delta_{13}$	$\Delta_{14}$	$\Delta_{15}$	$\Delta_{16}$	$\Delta_{17}$
$x_2$	$\Delta_{21}$	$\Delta_{22}$	$\Delta_{23}$	$\Delta_{24}$	$\Delta_{25}$	$\Delta_{26}$	$\Delta_{27}$
$x_3$	$\Delta_{31}$	$\Delta_{32}$	$\Delta_{33}$	$r_{3,4}^*$	$r_{3,5}^*$	$\Delta_{36}$	$\Delta_{37}$
$x_4$	$\Delta_{41}$	$\Delta_{42}$	$\Delta_{43}$	$r_{4,4}^*$	$r_{4,5}^*$	$\Delta_{46}$	$\Delta_{47}$
$x_5$	$\Delta_{51}$	$\Delta_{52}$	$\Delta_{53}$	$\Delta_{54}$	$\Delta_{55}$	$\Delta_{56}$	$\Delta_{57}$

$$r_{4,5}^* = \min(r_{3,5}^*, r_{4,4}^*, r_{3,4}^*) + \Delta_{4,5}$$

# Best Path: Bellman Recursion

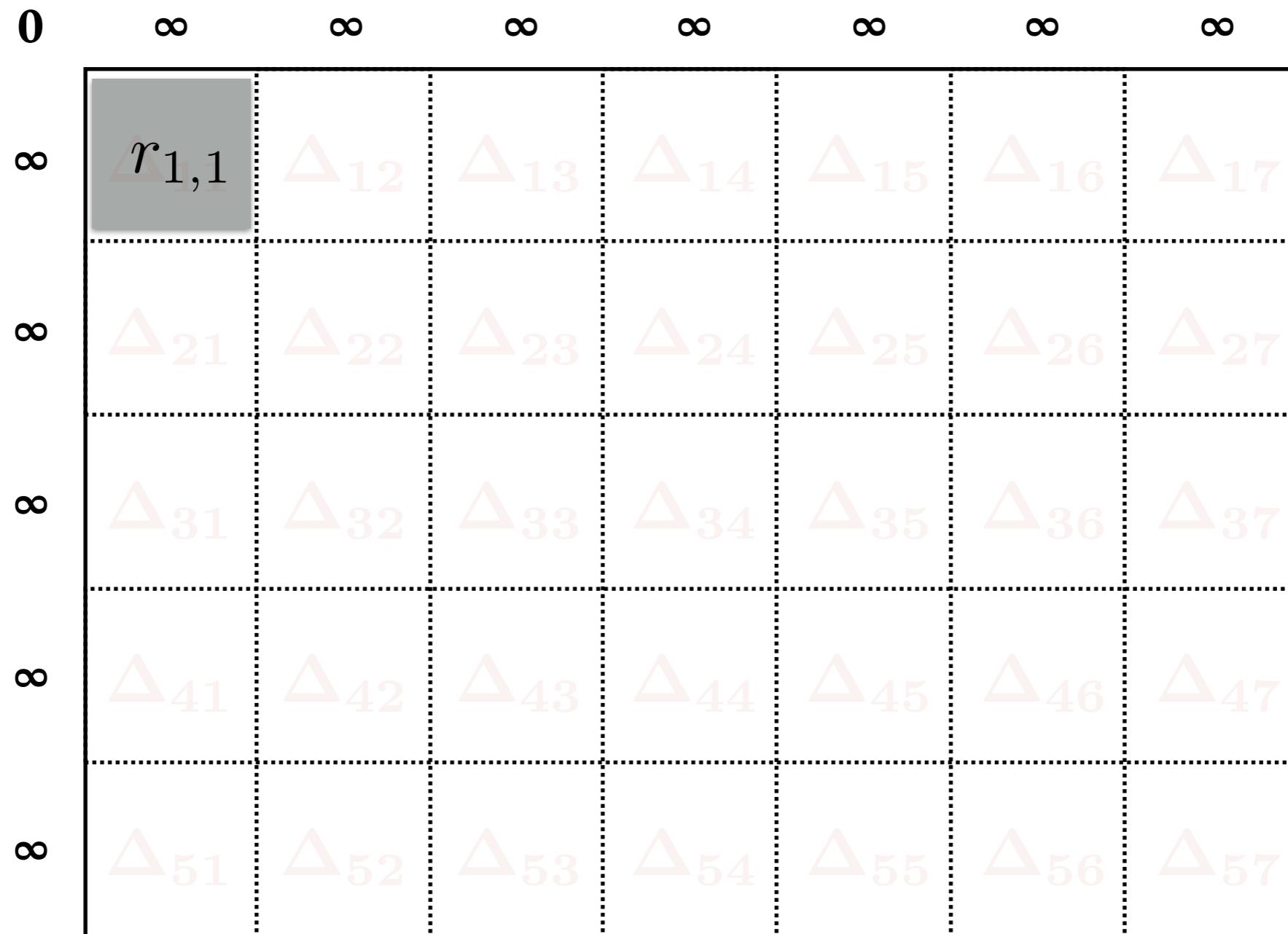


$$r_{1,1} = \Delta_{11}$$

$$r_{0,j} = r_{i,0} = \infty$$

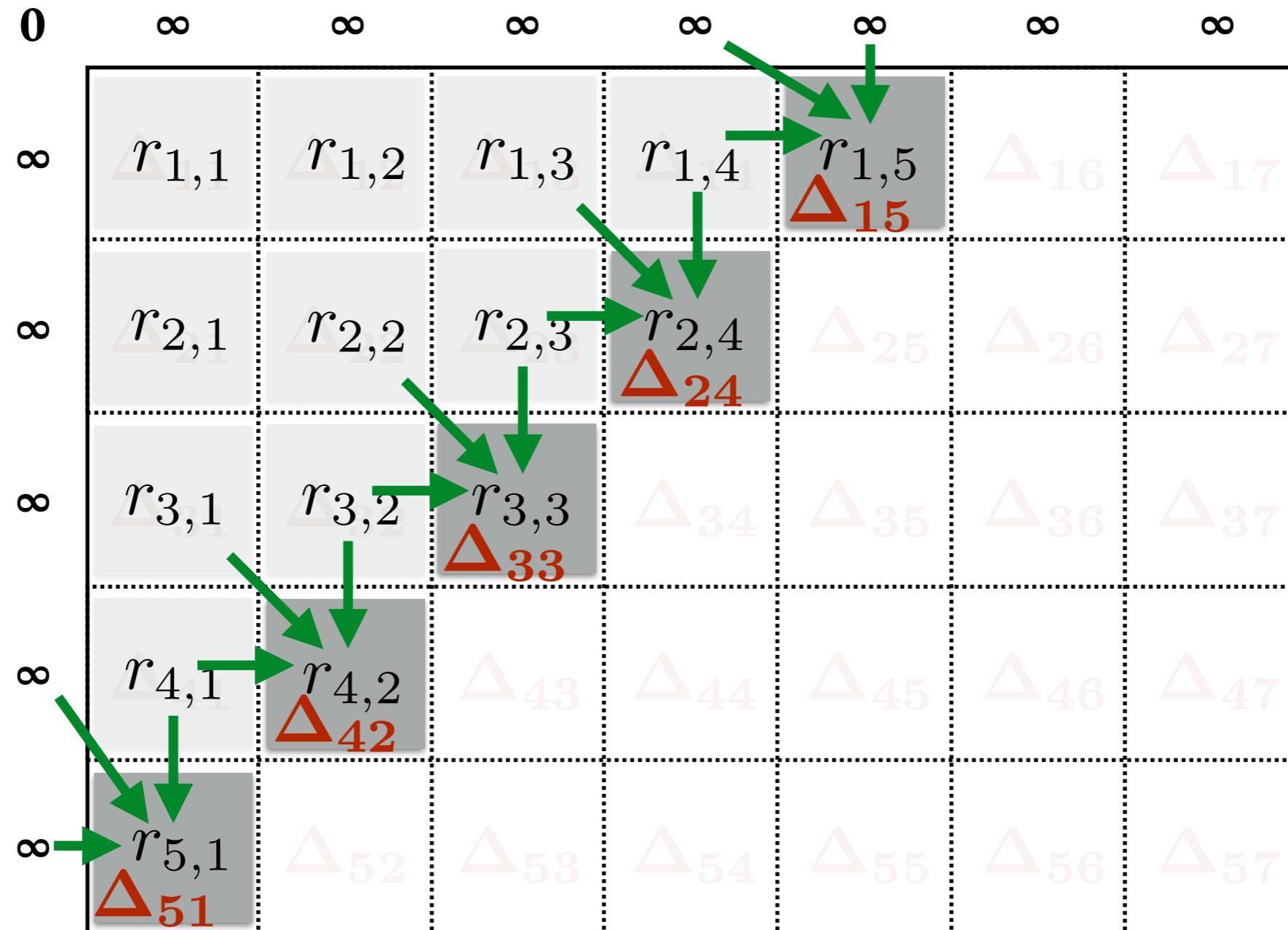
$$r_{0,0} = 0$$

# Best Path: Bellman Recursion



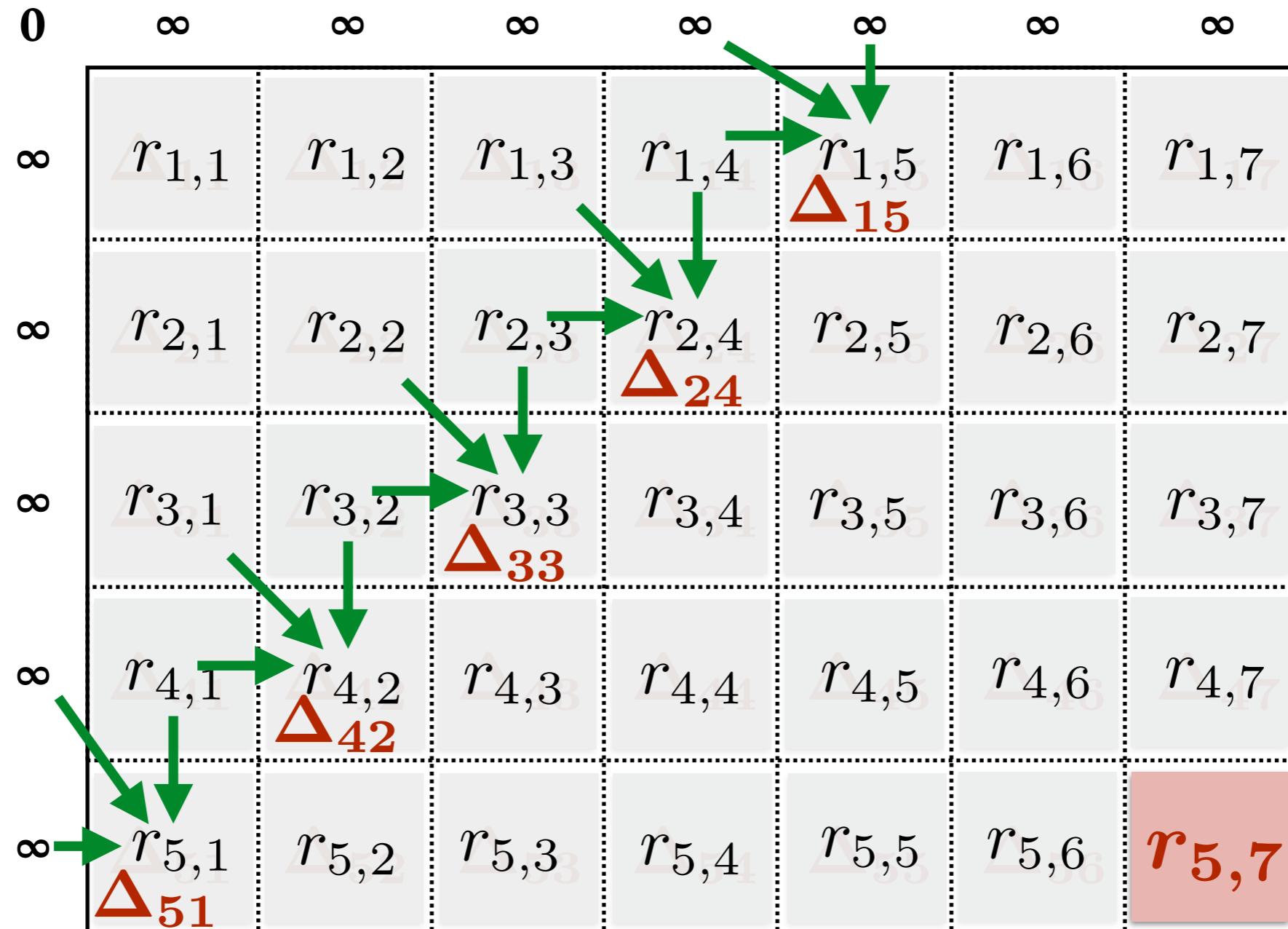
$$r_{i,j} = \min(r_{i-1,j-1}, r_{i-1,j}, r_{i,j-1}) + \Delta_{i,j}$$

# Best Path: Bellman Recursion



$$r_{i,j} = \min(r_{i-1,j-1}, r_{i-1,j}, r_{i,j-1}) + \Delta_{i,j}$$

# Best Path: Bellman Recursion



$$\text{dtw}_0(X, Y) = r_{n,m}$$

# Optimal Path

	0	$\infty$						
	$\infty$	1	0	0	0	0	0	0
	$\infty$	0	1	0	0	0	0	0
$A^*$	$\infty$	0	0	1	1	0	0	0
	$\infty$	0	0	0	0	1	0	0
	$\infty$	0	0	0	0	0	1	1

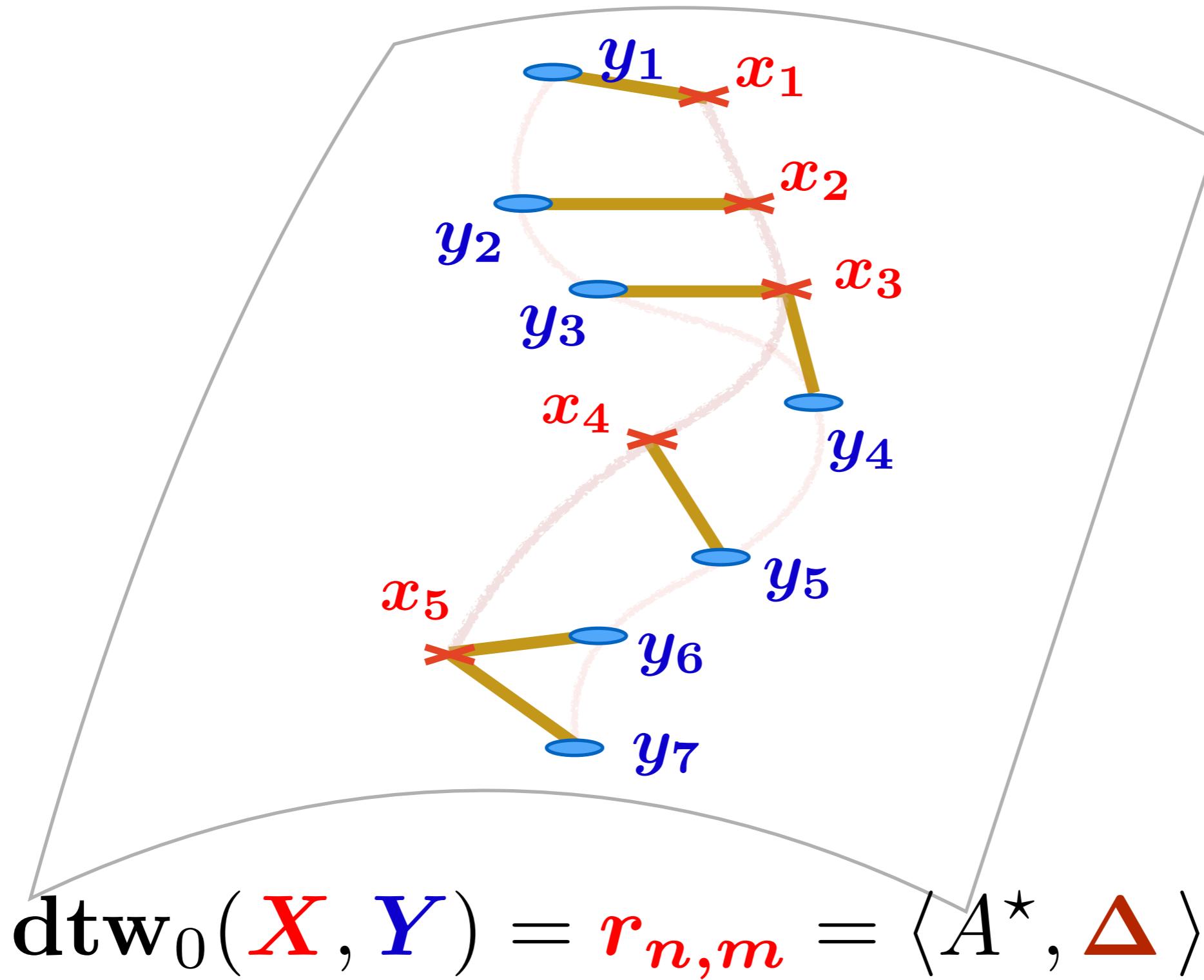
$$\text{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \textcolor{red}{r_{n,m}} = \langle A^*, \Delta \rangle$$

0. The DTW Geometry

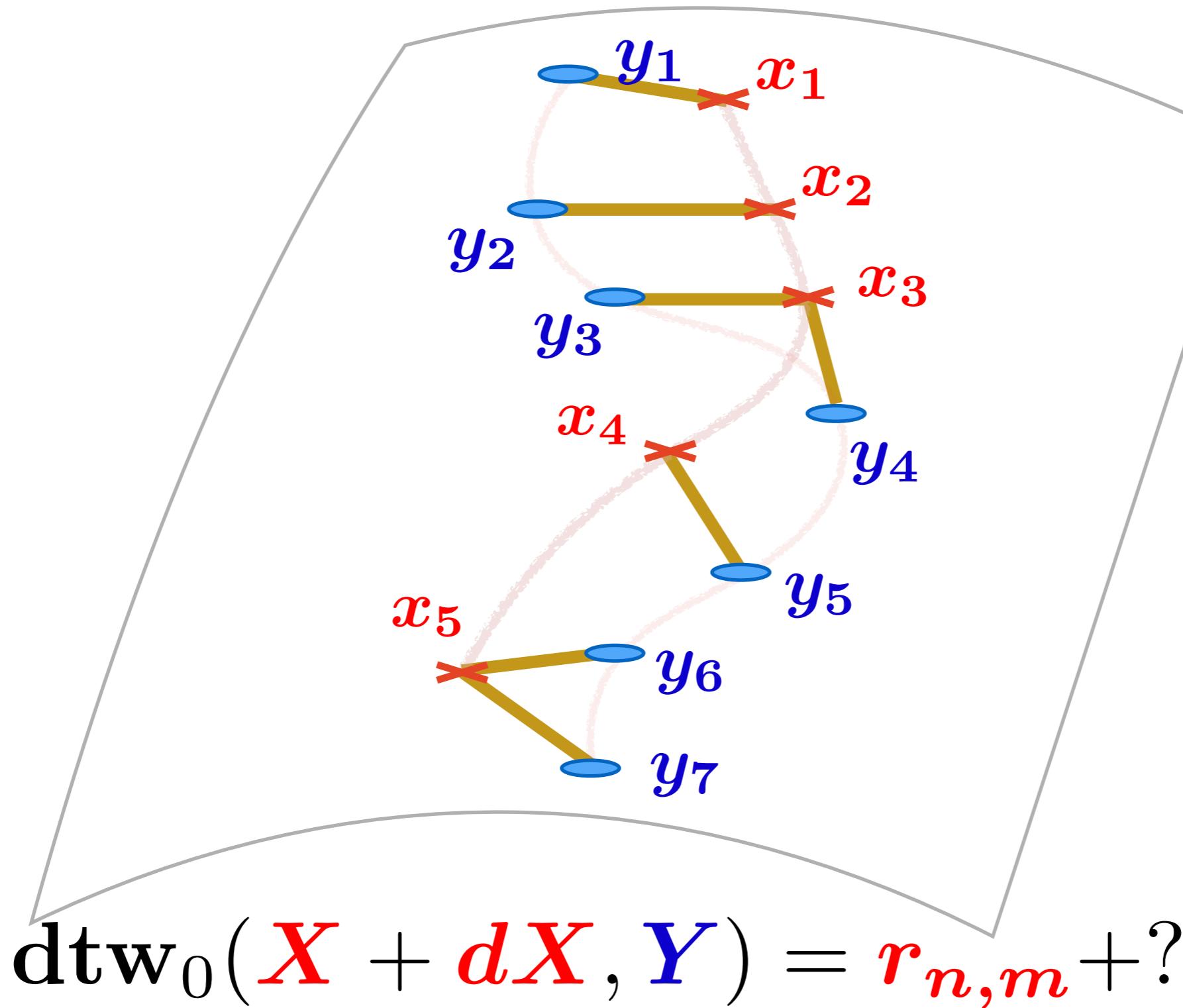
1. Soft-DTW

2. Soft-DTW as a Loss Function

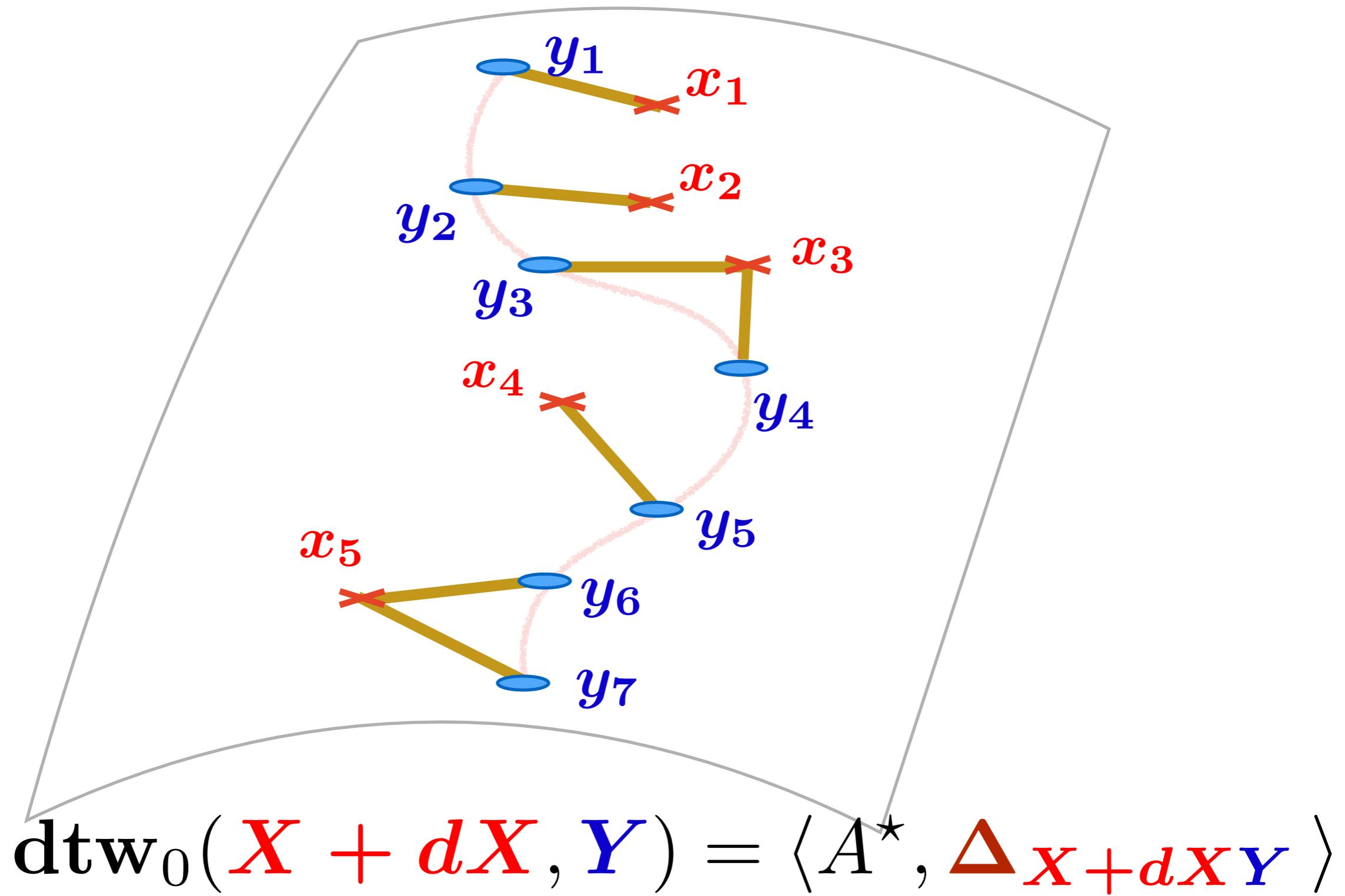
# DTW as a Loss: Differentiability?



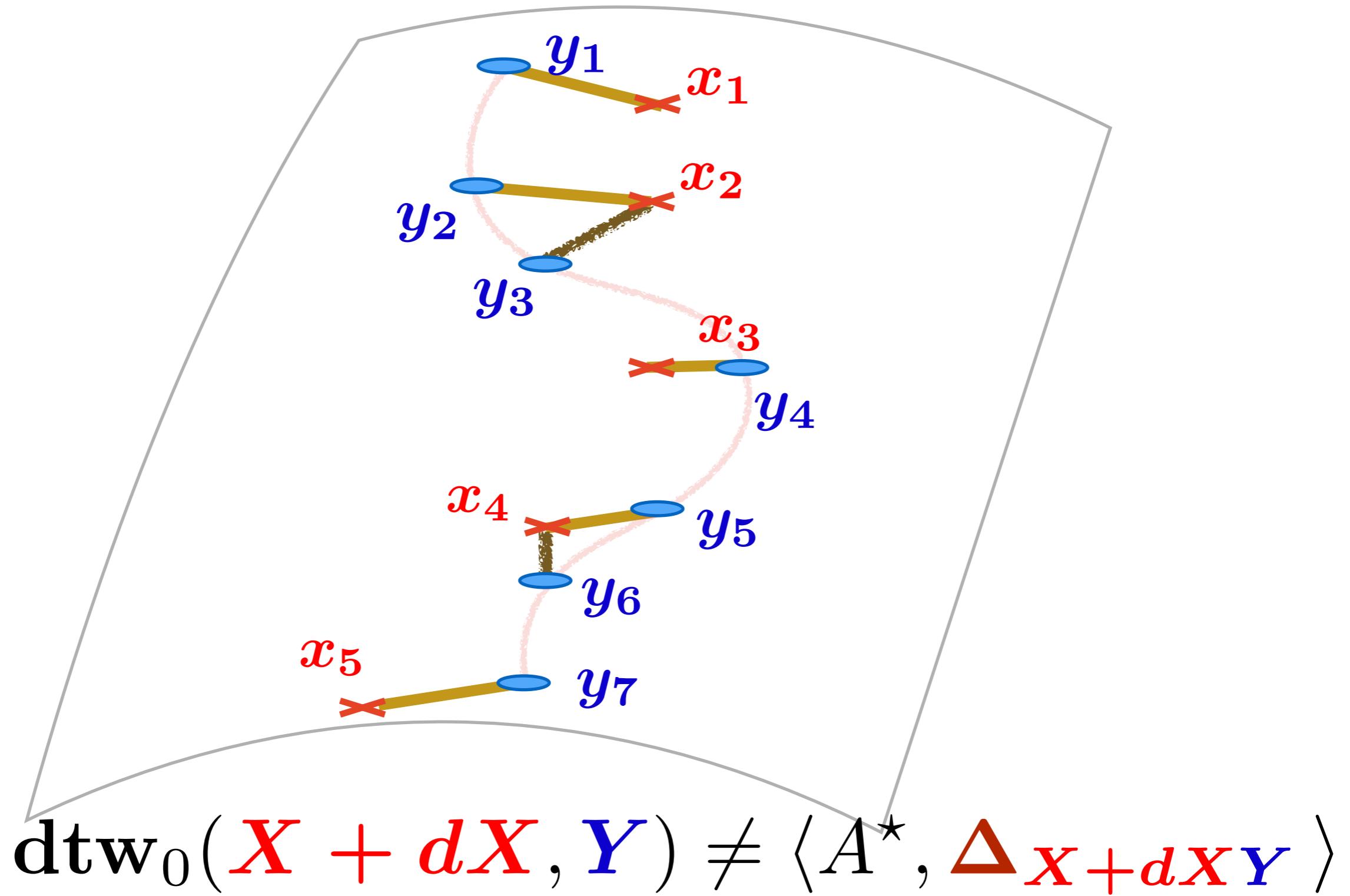
# DTW as a Loss: Differentiability?



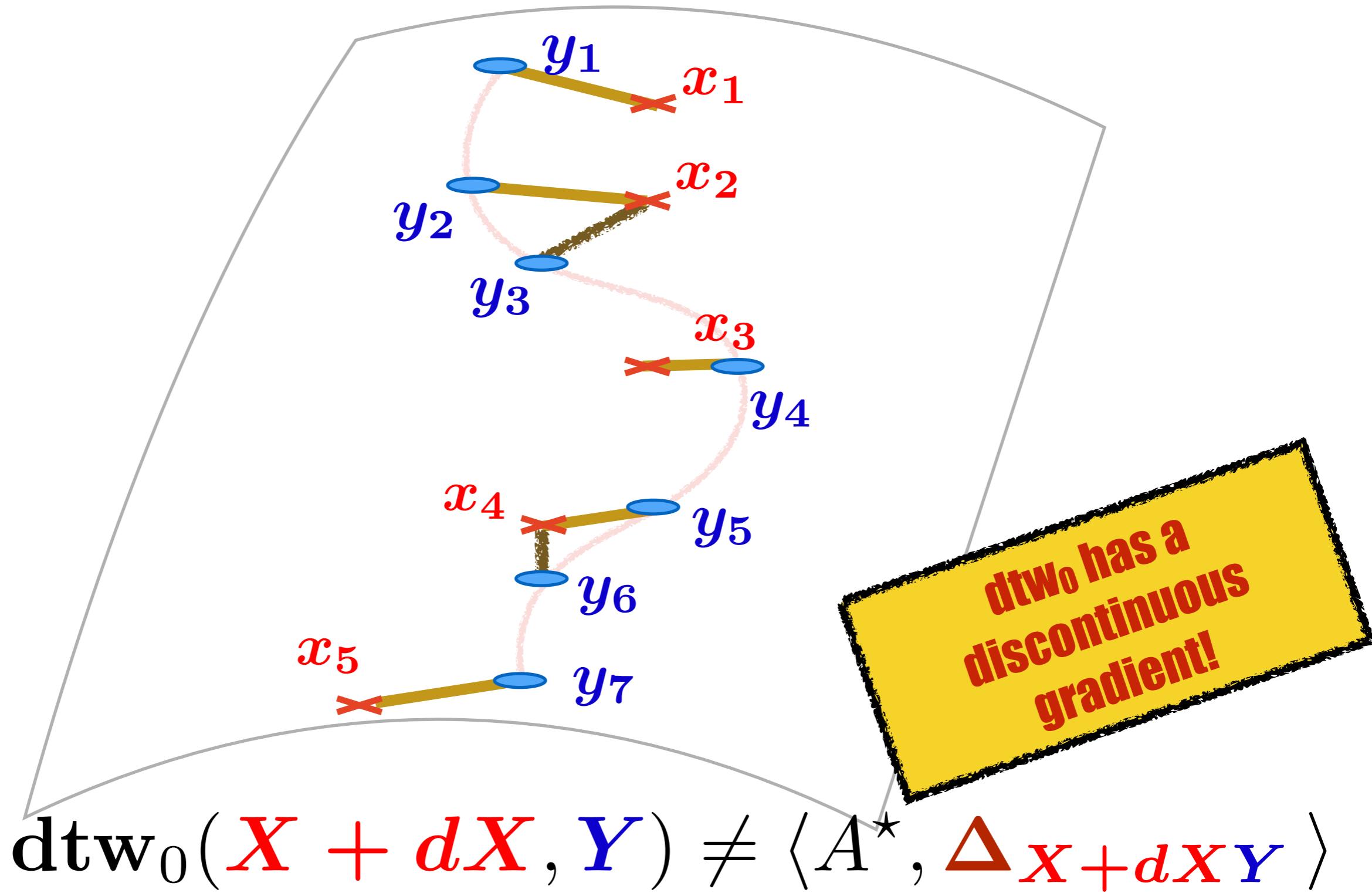
# DTW as a Loss: Differentiability?



# DTW as a Loss: Differentiability?



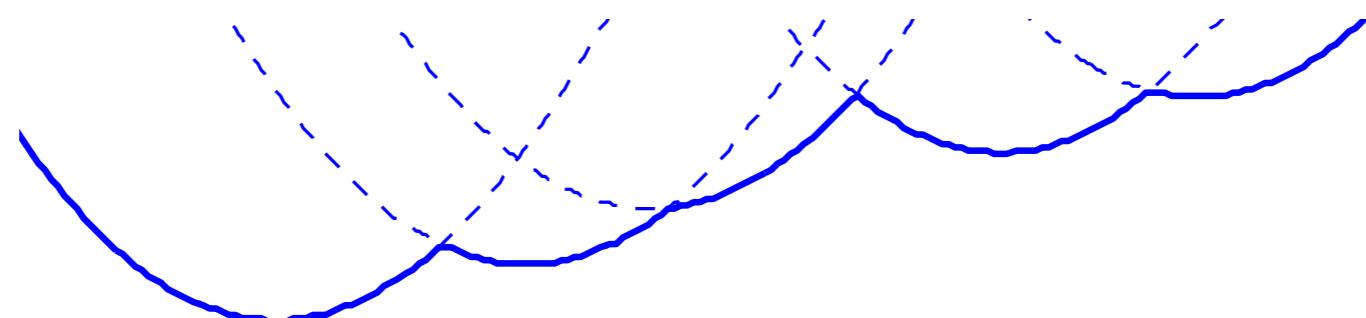
# DTW as a Loss: Differentiability?



# DTW as a Loss: Differentiability?

$$\text{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} \langle A, \Delta_{\textcolor{red}{X}\textcolor{blue}{Y}} \rangle$$

- $\text{dtw}_0$  is piecewise linear w.r.t  $\Delta$
- if  $\Delta_{ij} = \delta(\textcolor{red}{x}_i, \textcolor{blue}{y}_j) = \|\textcolor{red}{x}_i - \textcolor{blue}{y}_j\|^2$ ,  $\text{dtw}_0$  is piecewise quadratic w.r.t.  $\textcolor{red}{X}$ .



# DTW as a Loss: Differentiability?

$$\text{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} \langle A, \Delta_{\textcolor{red}{X}\textcolor{blue}{Y}} \rangle$$

$$\nabla_{\textcolor{red}{X}} \text{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \left( \frac{\partial \Delta_{\textcolor{red}{X}\textcolor{blue}{Y}}}{\partial \textcolor{red}{X}} \right)^T \quad \nabla_{\Delta} \min_{\mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} \langle \cdot, \Delta_{\textcolor{red}{X}\textcolor{blue}{Y}} \rangle$$

# DTW as a Loss: Differentiability?

$$\text{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} \langle A, \Delta_{\textcolor{red}{X}\textcolor{blue}{Y}} \rangle$$

$$\nabla_{\textcolor{red}{X}} \text{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \left( \frac{\partial \Delta_{\textcolor{red}{X}\textcolor{blue}{Y}}}{\partial \textcolor{red}{X}} \right)^T$$
$$\nabla_{\Delta} \min_{\mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} \langle \cdot, \Delta_{\textcolor{red}{X}\textcolor{blue}{Y}} \rangle$$

Jacobian matrix of  $\Delta$  w.r.t.  $\textcolor{red}{X}$

# DTW as a Loss: Differentiability?

$$\text{dtw}_0(\mathbf{X}, \mathbf{Y}) = \min_{A \in \mathcal{A}(\mathbf{n}, \mathbf{m})} \langle A, \Delta_{\mathbf{XY}} \rangle$$

$$\nabla_{\mathbf{X}} \text{dtw}_0(\mathbf{X}, \mathbf{Y}) = \left( \frac{\partial \Delta_{\mathbf{XY}}}{\partial \mathbf{X}} \right)^T$$
$$\nabla_{\Delta} \min_{\mathcal{A}(\mathbf{n}, \mathbf{m})} \langle \cdot, \Delta_{\mathbf{XY}} \rangle$$

Jacobian matrix of  $\Delta$  w.r.t.  $\mathbf{X}$

iff optimal solution  
is unique

$= A^*$

# DTW as a Loss: Differentiability?

$$\text{dtw}_0(\mathbf{X}, \mathbf{Y}) = \min_{A \in \mathcal{A}(\mathbf{n}, \mathbf{m})} \langle A, \Delta_{\mathbf{XY}} \rangle$$

$$\nabla_{\mathbf{X}} \text{dtw}_0(\mathbf{X}, \mathbf{Y}) = \left( \frac{\partial \Delta_{\mathbf{XY}}}{\partial \mathbf{X}} \right)^T$$
$$\nabla_{\Delta} \min_{\mathcal{A}(\mathbf{n}, \mathbf{m})} \langle \cdot, \Delta_{\mathbf{XY}} \rangle$$

Jacobian matrix of  $\Delta$  w.r.t.  $\mathbf{X}$

iff optimal solution  
is unique

$= A^*$

When  $A^*$  is not unique,  $\text{dtw}_0$  has a **discontinuous** gradient!

# Our proposal: smoothing the min

$$\text{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} \langle A, \Delta_{\textcolor{red}{X}\textcolor{blue}{Y}} \rangle$$

Problem: non-differentiability of min operator over finite family of values.

# Our proposal: smoothing the min

$$\text{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} \langle A, \Delta_{\textcolor{red}{X}\textcolor{blue}{Y}} \rangle$$

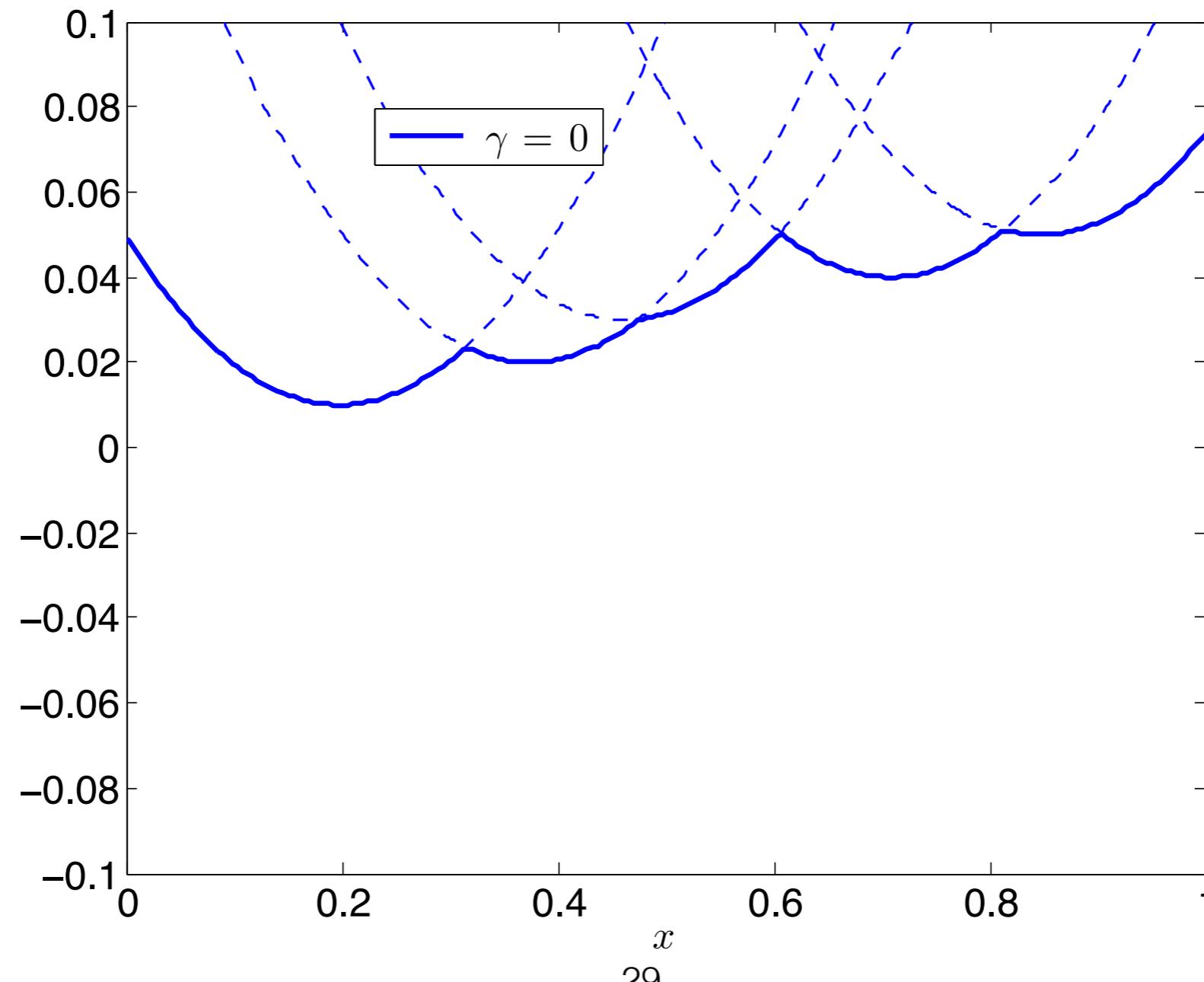
Problem: non-differentiability of min operator over finite family of values.

Fix: smoothed min operator

$$\min^\gamma(u_1, \dots, u_n) = \begin{cases} \min_{i \leq n} u_i, & \gamma = 0, \\ -\gamma \log \sum_{i=1}^n e^{-u_i/\gamma}, & \gamma > 0. \end{cases}$$

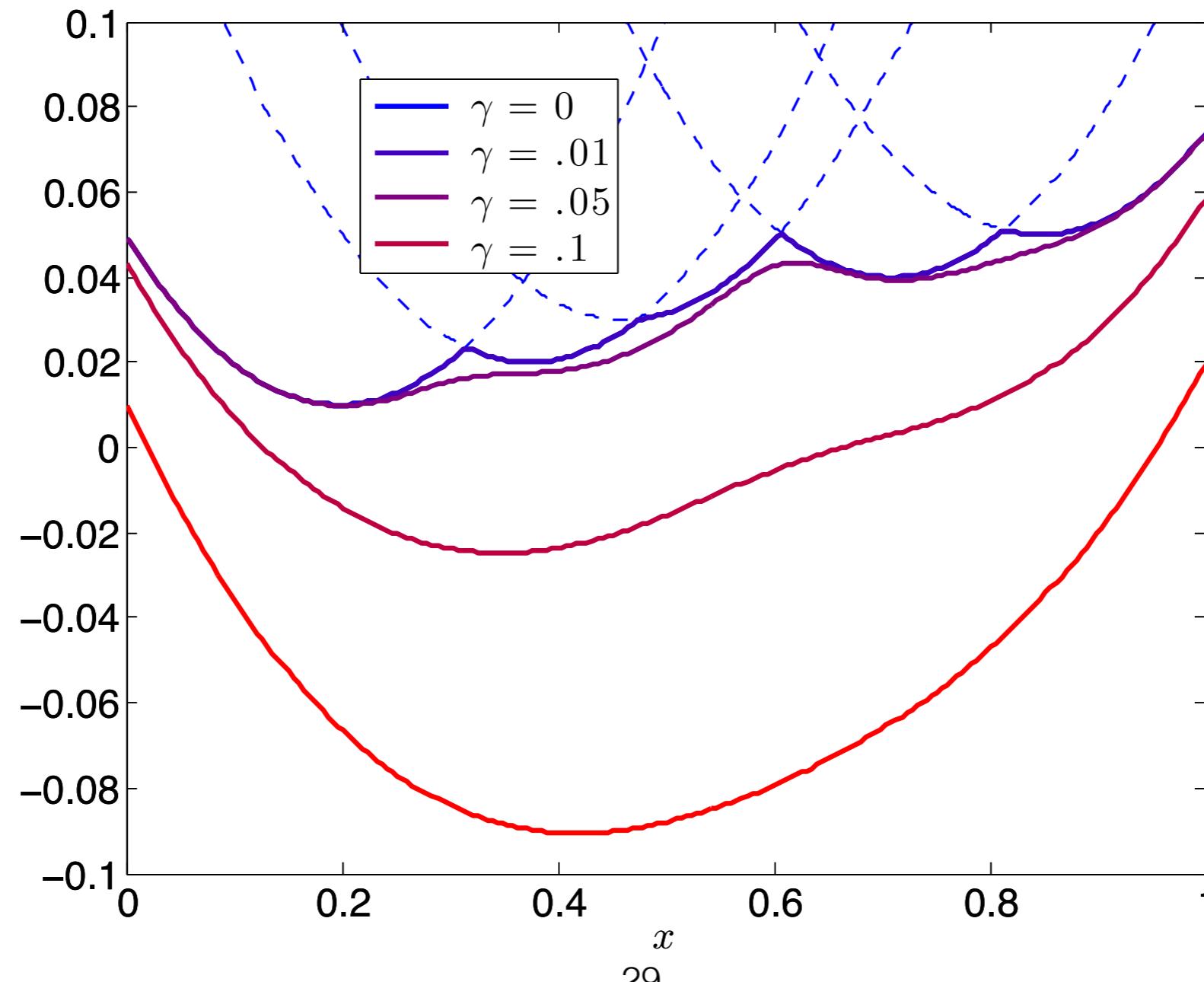
# *Example* softmin of quadratic functions

$$f(\boldsymbol{x}) = \min_{i=1,\dots,s}^{\gamma} a_i \boldsymbol{x}^2 + b_i \boldsymbol{x} + c_i$$



# *Example* softmin of quadratic functions

$$f(\mathbf{x}) = \min_{i=1,\dots,s}^{\gamma} a_i \mathbf{x}^2 + b_i \mathbf{x} + c_i$$



# Soft-DTW

$$\text{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} \langle A, \Delta_{\textcolor{red}{X}\textcolor{blue}{Y}} \rangle$$

# Soft-DTW

$$\text{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} \langle A, \Delta_{\textcolor{red}{X}\textcolor{blue}{Y}} \rangle$$

**Fix:** Replace  $\min$  by  $\min^\gamma$ ,  $\gamma > 0$

$$\text{dtw}_{\gamma}(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})}^{\gamma} \langle A, \Delta_{\textcolor{red}{X}\textcolor{blue}{Y}} \rangle$$

# Soft-DTW

$$\text{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} \langle A, \Delta_{\textcolor{red}{X} \textcolor{blue}{Y}} \rangle$$

**Fix:** Replace  $\min$  by  $\min^\gamma$ ,  $\gamma > 0$

$$\text{dtw}_\gamma(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})}^{\gamma} \langle A, \Delta_{\textcolor{red}{X} \textcolor{blue}{Y}} \rangle$$

$$\text{dtw}_\gamma(\textcolor{red}{X}, \textcolor{blue}{Y}) = -\gamma \log \sum_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} e^{-\frac{\langle A, \Delta_{\textcolor{red}{X} \textcolor{blue}{Y}} \rangle}{\gamma}}$$

# Relation to Global Alignment kernels

$$k_{\text{GA}} := \sum_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} e^{-\frac{\langle A, \Delta_{\textcolor{red}{X} \textcolor{blue}{Y}} \rangle}{\gamma}}$$



A positive semi-definite **kernel** between time series

# Relation to Global Alignment kernels

$$k_{\text{GA}} := \sum_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} e^{-\frac{\langle A, \Delta_{\textcolor{red}{X} \textcolor{blue}{Y}} \rangle}{\gamma}}$$



A positive semi-definite **kernel** between time series

$$\text{dtw}_{\gamma}(\textcolor{red}{X}, \textcolor{blue}{Y}) = -\gamma \log k_{\text{GA}}$$

Computing soft-DTW is equivalent to computing  $k_{\text{GA}}$  in **log domain**

# Recursive Computation

$$\text{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} \langle A, \Delta_{\textcolor{red}{X}\textcolor{blue}{Y}} \rangle$$

$$r_{i,j} = \min(r_{i-1,j-1}, r_{i-1,j}, r_{i,j-1}) + \Delta_{i,j}$$

# Recursive Computation

$$\text{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} \langle A, \Delta_{\textcolor{red}{X}\textcolor{blue}{Y}} \rangle$$

$$r_{i,j} = \min(r_{i-1,j-1}, r_{i-1,j}, r_{i,j-1}) + \Delta_{i,j}$$

$$\text{dtw}_{\gamma}(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})}^{\gamma} \langle A, \Delta_{\textcolor{red}{X}\textcolor{blue}{Y}} \rangle$$

$$r_{i,j} = \min^{\gamma}(r_{i-1,j-1}, r_{i-1,j}, r_{i,j-1}) + \Delta_{i,j}$$

# Recursive Computation

$$\text{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} \langle A, \Delta_{\textcolor{red}{X}\textcolor{blue}{Y}} \rangle$$

$$r_{i,j} = \min(r_{i-1,j-1}, r_{i-1,j}, r_{i,j-1}) + \Delta_{i,j}$$

$$\text{dtw}_{\gamma}(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})}^{\gamma} \langle A, \Delta_{\textcolor{red}{X}\textcolor{blue}{Y}} \rangle$$

$$r_{i,j} = \min^{\gamma}(r_{i-1,j-1}, r_{i-1,j}, r_{i,j-1}) + \Delta_{i,j}$$

Simply replace min operator!

# Recursive Computation

$$\text{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} \langle A, \Delta_{\textcolor{red}{X}\textcolor{blue}{Y}} \rangle$$

$$r_{i,j} = \min(r_{i-1,j-1}, r_{i-1,j}, r_{i,j-1}) + \Delta_{i,j}$$

$$\text{dtw}_{\gamma}(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})}^{\gamma} \langle A, \Delta_{\textcolor{red}{X}\textcolor{blue}{Y}} \rangle$$

$$r_{i,j} = \min^{\gamma}(r_{i-1,j-1}, r_{i-1,j}, r_{i,j-1}) + \Delta_{i,j}$$

Simply replace min operator!

Stable: recursion in log domain!

# Differentiation

$$\text{dtw}_{\gamma}(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})}^{\gamma} \langle A, \Delta_{\textcolor{red}{X} \textcolor{blue}{Y}} \rangle$$

$$\nabla_X \text{dtw}_0(\textcolor{red}{X}, \textcolor{blue}{Y}) = \left( \frac{\partial \Delta(\textcolor{red}{X}, \textcolor{blue}{Y})}{\partial \textcolor{red}{X}} \right)^T A^\star$$

# Differentiation

$$\text{dtw}_{\gamma}(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})}^{\gamma} \langle A, \Delta_{\textcolor{red}{X} \textcolor{blue}{Y}} \rangle$$

$$\nabla_X \text{dtw}_{\gamma}(\textcolor{red}{X}, \textcolor{blue}{Y}) = \left( \frac{\partial \Delta(\textcolor{red}{X}, \textcolor{blue}{Y})}{\partial \textcolor{red}{X}} \right)^T \mathbb{E}_{\gamma}[A]$$

# Differentiation

$$\text{dtw}_{\gamma}(X, Y) = \min_{A \in \mathcal{A}(n, m)}^{\gamma} \langle A, \Delta_{XY} \rangle$$

$$\nabla_X \text{dtw}_{\gamma}(X, Y) = \left( \frac{\partial \Delta(X, Y)}{\partial X} \right)^T \mathbb{E}_{\gamma}[A]$$

$$\mathbb{E}_{\gamma}[A] := \frac{\sum_{A \in \mathcal{A}(n, m)} A e^{-\frac{\langle A, \Delta_{XY} \rangle}{\gamma}}}{\sum_{A \in \mathcal{A}(n, m)} e^{-\frac{\langle A, \Delta_{XY} \rangle}{\gamma}}}$$

**Expectation of Path  
under Gibbs  
distribution**

# Differentiation

$$\text{dtw}_{\gamma}(\textcolor{red}{X}, \textcolor{blue}{Y}) = \min_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})}^{\gamma} \langle A, \Delta_{\textcolor{red}{X} \textcolor{blue}{Y}} \rangle$$

$$\nabla_X \text{dtw}_{\gamma}(\textcolor{red}{X}, \textcolor{blue}{Y}) = \left( \frac{\partial \Delta(\textcolor{red}{X}, \textcolor{blue}{Y})}{\partial \textcolor{red}{X}} \right)^T \mathbb{E}_{\gamma}[A]$$

$\nabla_{\Delta} \text{dtw}_{\gamma}(\textcolor{red}{X}, \textcolor{blue}{Y})$

$$\mathbb{E}_{\gamma}[A] := \frac{\sum_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} A e^{-\frac{\langle A, \Delta_{\textcolor{red}{X} \textcolor{blue}{Y}} \rangle}{\gamma}}}{\sum_{A \in \mathcal{A}(\textcolor{red}{n}, \textcolor{blue}{m})} e^{-\frac{\langle A, \Delta_{\textcolor{red}{X} \textcolor{blue}{Y}} \rangle}{\gamma}}}$$

**Expectation of Path  
under Gibbs  
distribution**

# Computing the expectation $E_{\gamma}[A]$

$$E_{\gamma}[A] := \frac{\sum_{A \in \mathcal{A}(\mathbf{n}, \mathbf{m})} A e^{-\frac{\langle A, \Delta_{\mathbf{X} \mathbf{Y}} \rangle}{\gamma}}}{\sum_{A \in \mathcal{A}(\mathbf{n}, \mathbf{m})} e^{-\frac{\langle A, \Delta_{\mathbf{X} \mathbf{Y}} \rangle}{\gamma}}}$$

Naive computation  
is intractable

# Computing the expectation $E_{\gamma}[A]$

$$E_{\gamma}[A] := \frac{\sum_{A \in \mathcal{A}(\mathbf{n}, \mathbf{m})} A e^{-\frac{\langle A, \Delta_{\mathbf{X} \mathbf{Y}} \rangle}{\gamma}}}{\sum_{A \in \mathcal{A}(\mathbf{n}, \mathbf{m})} e^{-\frac{\langle A, \Delta_{\mathbf{X} \mathbf{Y}} \rangle}{\gamma}}}$$

Naive computation  
is intractable

$$= \frac{\sum_{A \in \mathcal{A}(\mathbf{n}, \mathbf{m})} A e^{-\frac{\langle A, \Delta_{\mathbf{X} \mathbf{Y}} \rangle}{\gamma}}}{k_{GA}}$$

$k_{GA}$  is the  
normalization constant  
(a.k.a. partition function)!

# Computing the expectation $E_{\gamma}[A]$

$$E_{\gamma}[A] := \frac{\sum_{A \in \mathcal{A}(\mathbf{n}, \mathbf{m})} A e^{-\frac{\langle A, \Delta_{\mathbf{X} \mathbf{Y}} \rangle}{\gamma}}}{\sum_{A \in \mathcal{A}(\mathbf{n}, \mathbf{m})} e^{-\frac{\langle A, \Delta_{\mathbf{X} \mathbf{Y}} \rangle}{\gamma}}}$$

Naive computation  
is intractable

$$= \frac{\sum_{A \in \mathcal{A}(\mathbf{n}, \mathbf{m})} A e^{-\frac{\langle A, \Delta_{\mathbf{X} \mathbf{Y}} \rangle}{\gamma}}}{k_{GA}}$$

$k_{GA}$  is the  
normalization constant  
(a.k.a. partition function)!

$$= \nabla_{\Delta} -\gamma \log k_{GA}$$

$E_{\gamma}[A]$  is the gradient  
of the log partition

# Computing the expectation $E_{\gamma}[A]$

$$E_{\gamma}[A] := \frac{\sum_{A \in \mathcal{A}(\mathbf{n}, \mathbf{m})} Ae^{-\frac{\langle A, \Delta_{\mathbf{X} \mathbf{Y}} \rangle}{\gamma}}}{\sum_{A \in \mathcal{A}(\mathbf{n}, \mathbf{m})} e^{-\frac{\langle A, \Delta_{\mathbf{X} \mathbf{Y}} \rangle}{\gamma}}}$$

Naive computation  
is intractable

$$= \frac{\sum_{A \in \mathcal{A}(\mathbf{n}, \mathbf{m})} Ae^{-\frac{\langle A, \Delta_{\mathbf{X} \mathbf{Y}} \rangle}{\gamma}}}{k_{GA}}$$

$k_{GA}$  is the  
normalization constant  
(a.k.a. partition function)!

$$= \nabla_{\Delta} -\gamma \log k_{GA}$$



# Computing the expectation $E_{\gamma}[A]$

---

To summarize, we want to compute:

$$E_{\gamma}[A] = \nabla_{\Delta} -\gamma \log k_G A$$

# Computing the expectation $E_{\gamma}[A]$

---

To summarize, we want to compute:

$$E_{\gamma}[A] = \nabla_{\Delta} -\gamma \log k_{GA} = \nabla_{\Delta} dtw_{\gamma}$$

# Computing the expectation $E_{\gamma}[A]$

---

To summarize, we want to compute:

$$E_{\gamma}[A] = \nabla_{\Delta} -\gamma \log k_{GA} = \nabla_{\Delta} dtw_{\gamma}$$

$E_{\gamma}[A]$  can be computed by **backpropagation** in the same  $O(nm)$  cost as  $dtw_{\gamma}$

# Computing the expectation $E_{\gamma}[A]$

To summarize, we want to compute:

$$E_{\gamma}[A] = \nabla_{\Delta} -\gamma \log k_{GA} = \nabla_{\Delta} dtw_{\gamma}$$

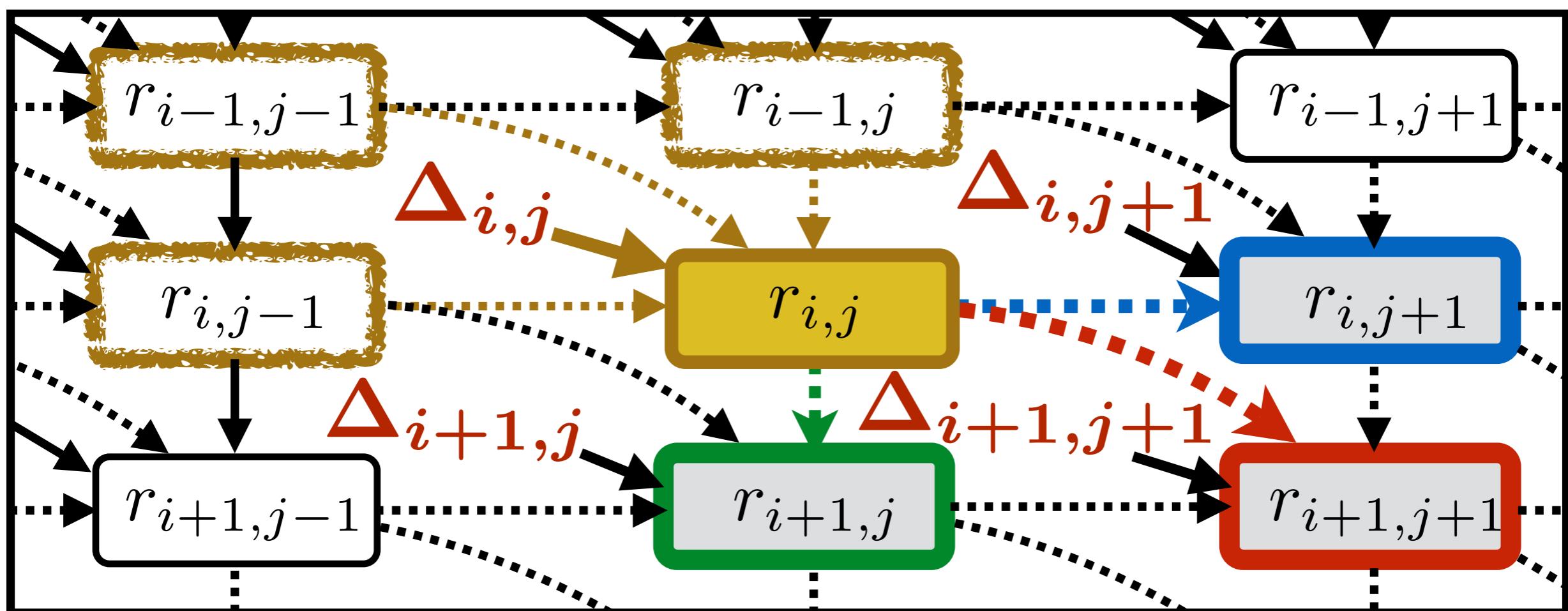
$E_{\gamma}[A]$  can be computed by **backpropagation** in the same  $O(nm)$  cost as  $dtw_{\gamma}$

We derive a backward recursion **without resorting to autodiff**

Faster and more numerically stable

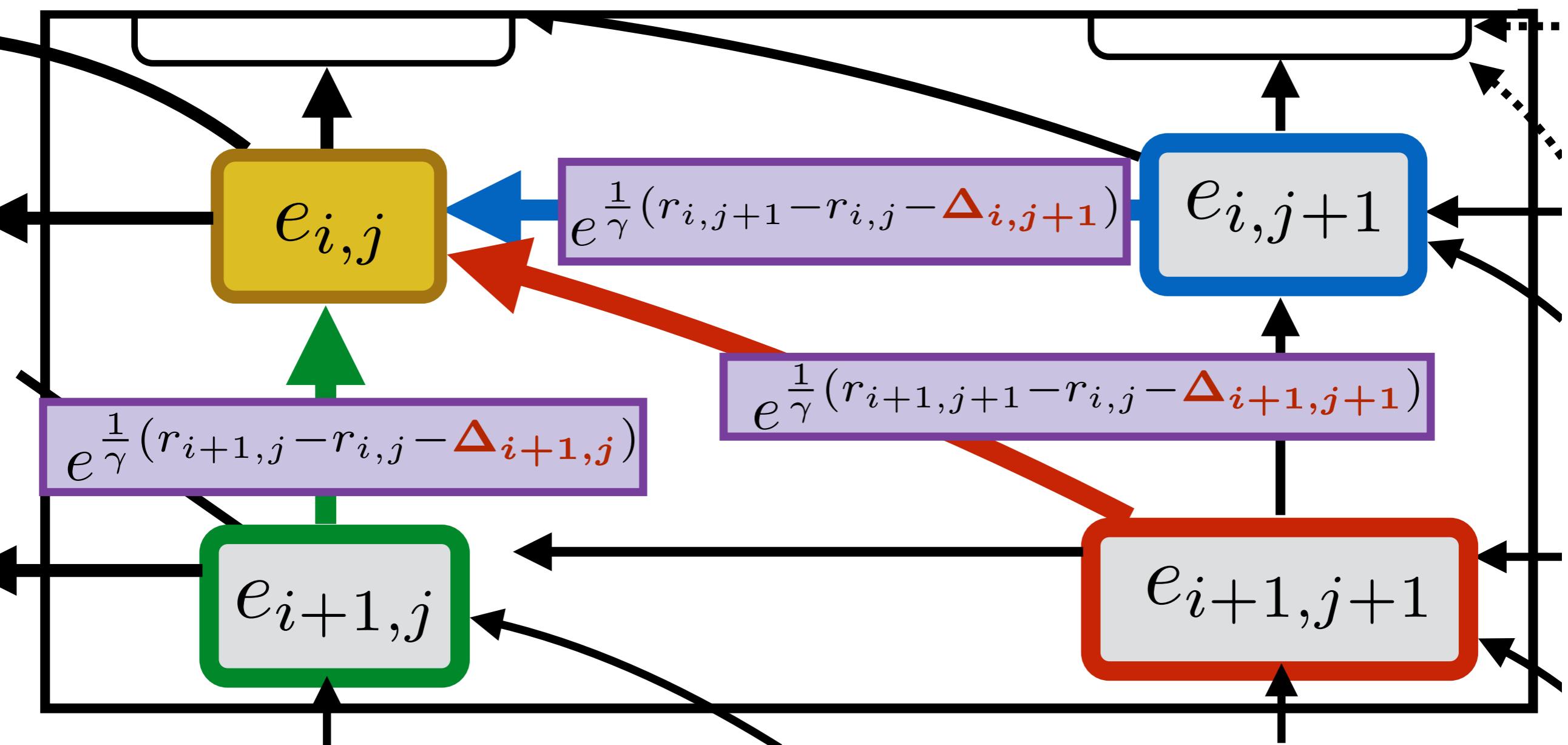
# Forward Pass

Bellman's recursion has the following computational graph

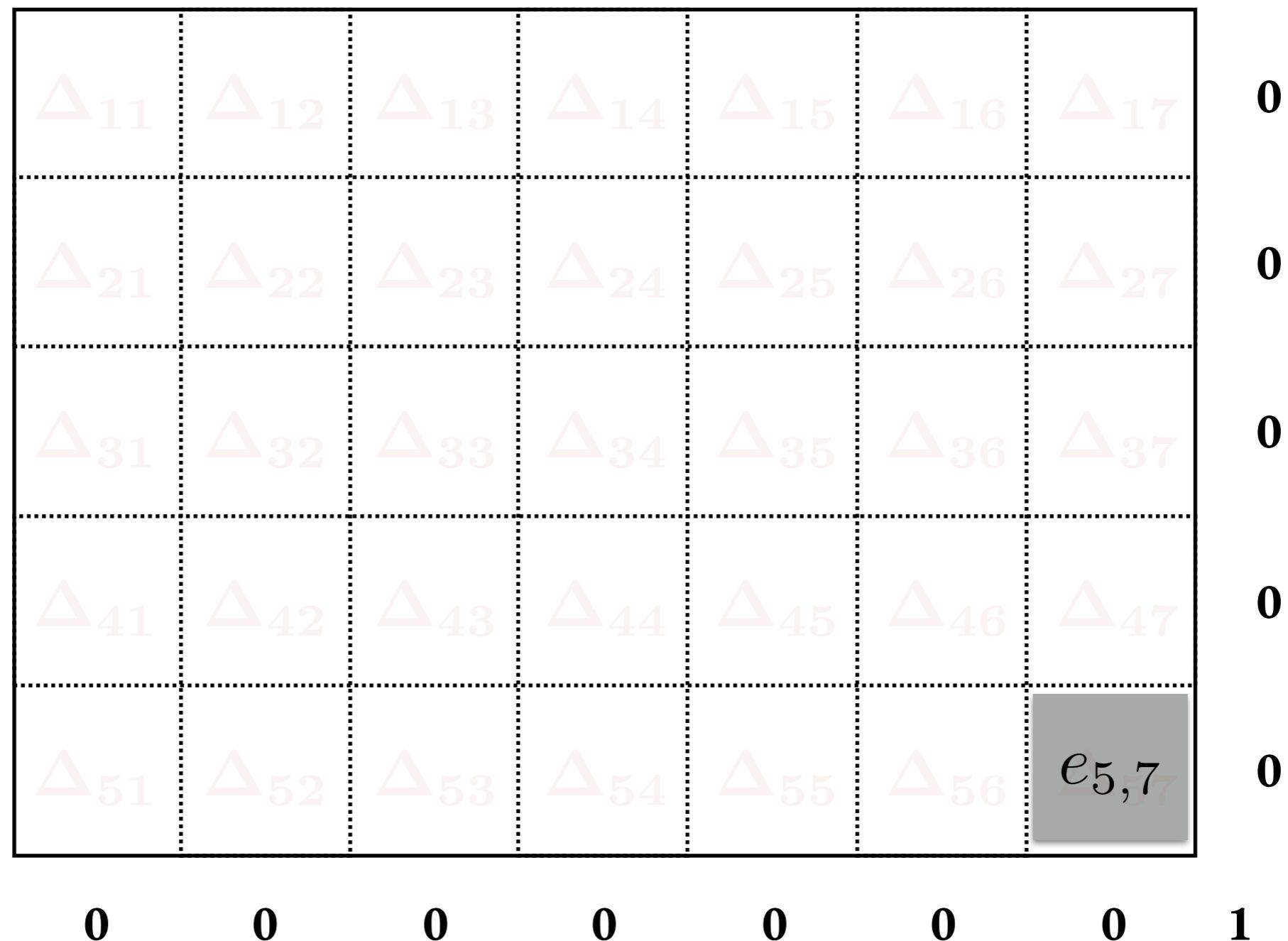


# Backward Pass

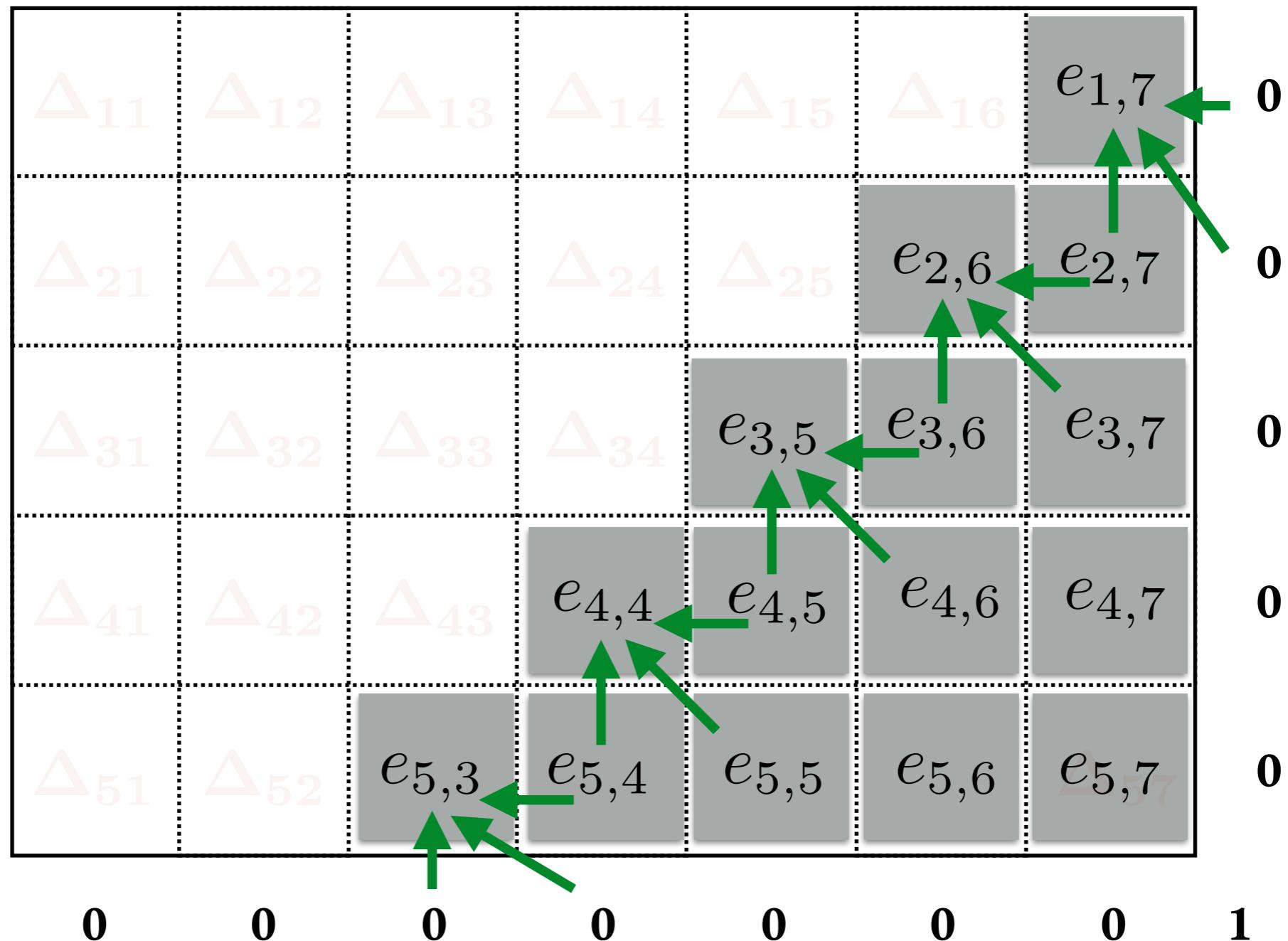
with a few simplifications, the backward pass boils down to the following updates



# Backward Recursion



# Backward Recursion



# Backward Recursion

$e_{1,1}$	$e_{1,2}$	$e_{1,3}$	$e_{1,4}$	$e_{1,5}$	$e_{1,6}$	$e_{1,7}$	<b>0</b>
$e_{2,1}$	$e_{2,2}$	$e_{2,3}$	$e_{2,4}$	$e_{2,5}$	$e_{2,6}$	$e_{2,7}$	<b>0</b>
$e_{3,1}$	$e_{3,2}$	$e_{3,3}$	$e_{3,4}$	$e_{3,5}$	$e_{3,6}$	$e_{3,7}$	<b>0</b>
$e_{4,1}$	$e_{4,2}$	$e_{4,3}$	$e_{4,4}$	$e_{4,5}$	$e_{4,6}$	$e_{4,7}$	<b>0</b>
$e_{5,1}$	$e_{5,2}$	$e_{5,3}$	$e_{5,4}$	$e_{5,5}$	$e_{5,6}$	$e_{5,7}$	<b>0</b>
0	0	0	0	0	0	0	1

$$E_{\gamma}[A] = [e]_{ij}$$

# Backward Recursion

$$a = e^{\frac{1}{\gamma}}(r_{i+1,j} - r_{i,j} - \Delta_{i+1,j})$$

$$b = e^{\frac{1}{\gamma}}(r_{i,j+1} - r_{i,j} - \Delta_{i,j+1})$$

$$c = e^{\frac{1}{\gamma}}(r_{i+1,j+1} - r_{i,j} - \Delta_{i+1,j+1})$$

$$e_{i,j} = e_{i+1,j} \cdot a + e_{i,j+1} \cdot b + e_{i+1,j+1} \cdot c$$

$$\nabla_X \text{dtw}_\gamma(\textcolor{red}{X}, \textcolor{blue}{Y}) = \left( \frac{\partial \Delta(\textcolor{red}{X}, \textcolor{blue}{Y})}{\partial \textcolor{red}{X}} \right)^T E$$

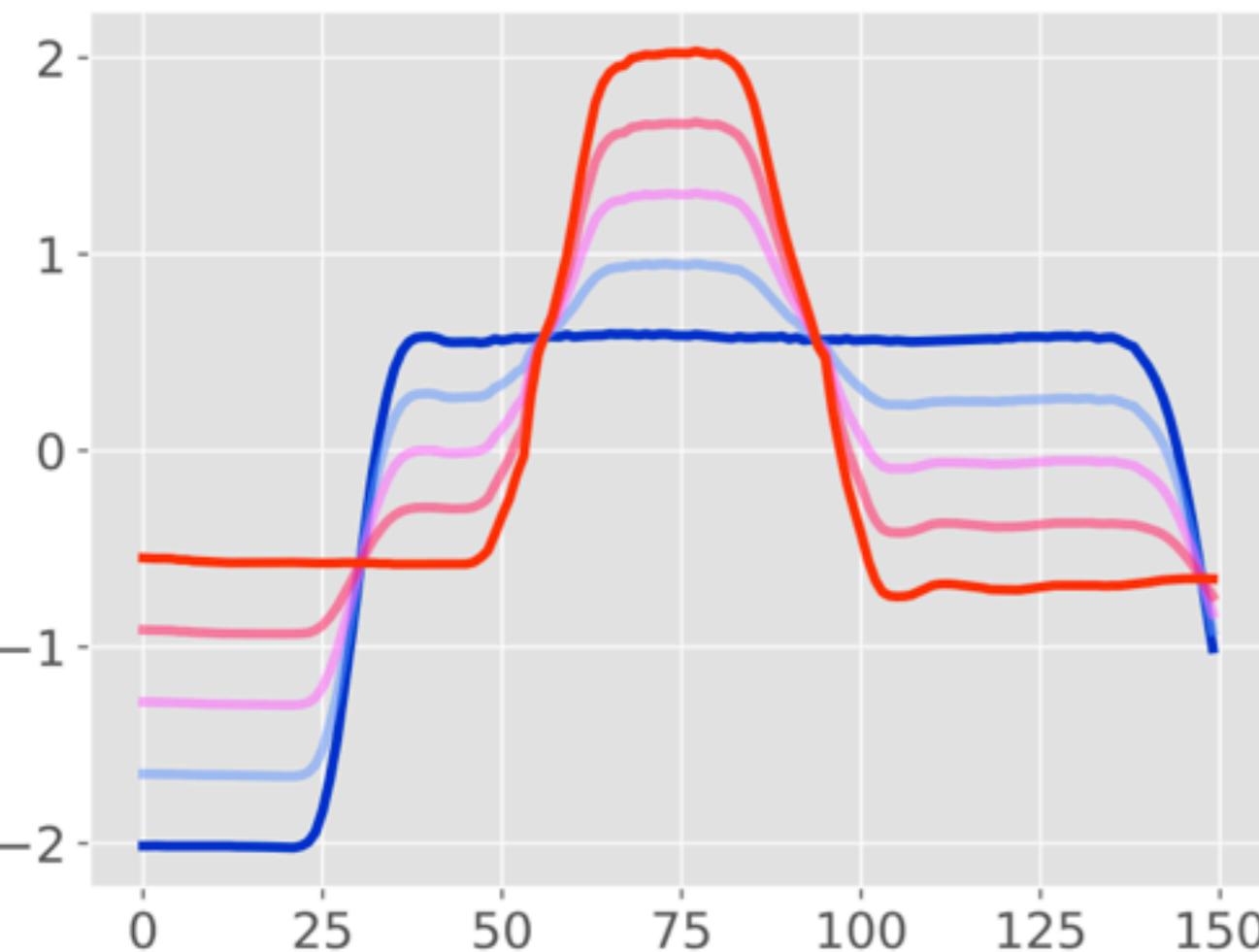
## 0. The DTW Geometry

### 1. Soft-DTW

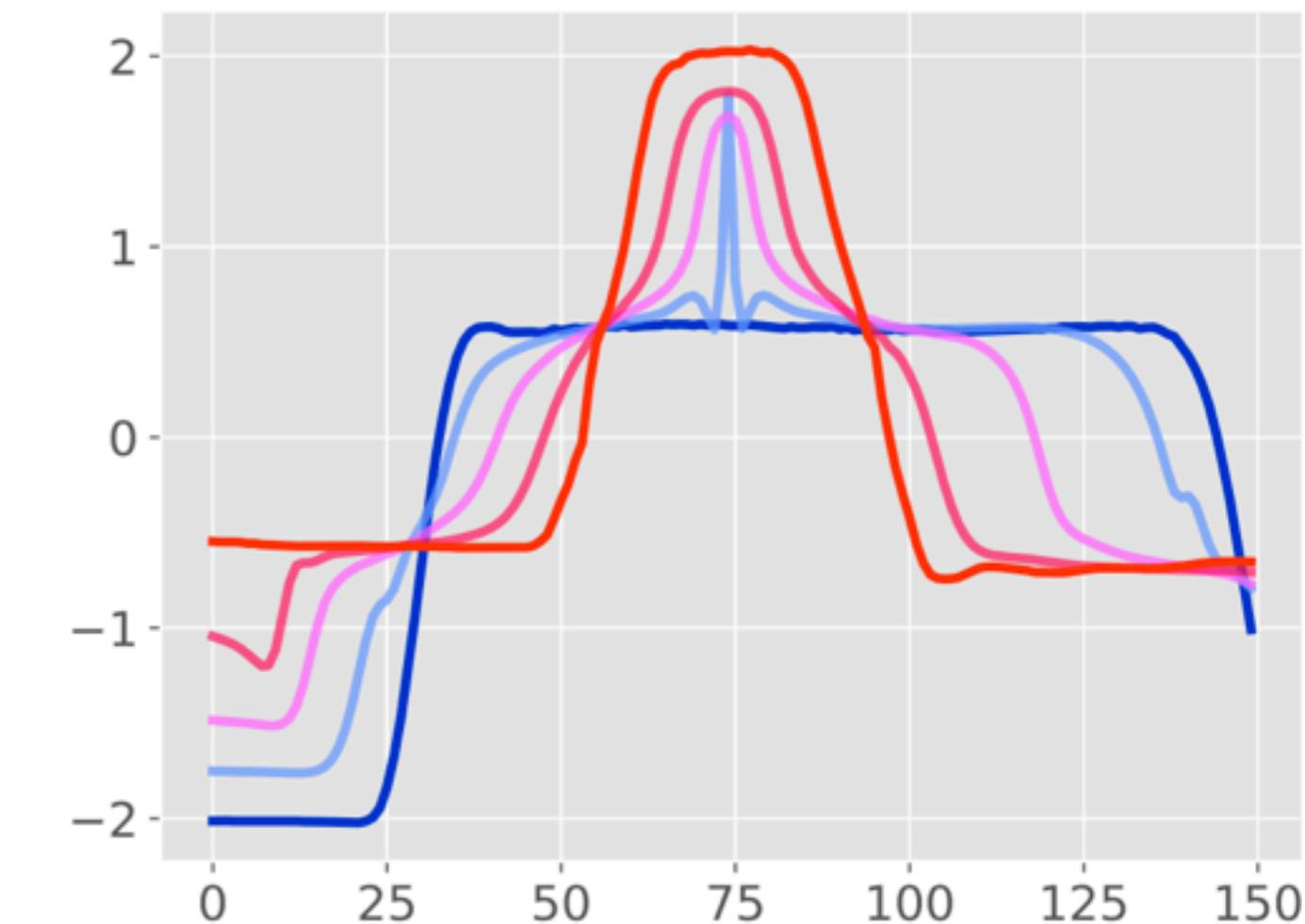
### 2. Soft-DTW as a Loss Function

# Interpolation Between 2 Time Series

$$\min_{\mathbf{X}} [\lambda \operatorname{dtw}_{\gamma}(\mathbf{X}, \mathbf{Y}_1) + (1 - \lambda) \operatorname{dtw}_{\gamma}(\mathbf{X}, \mathbf{Y}_2)]$$



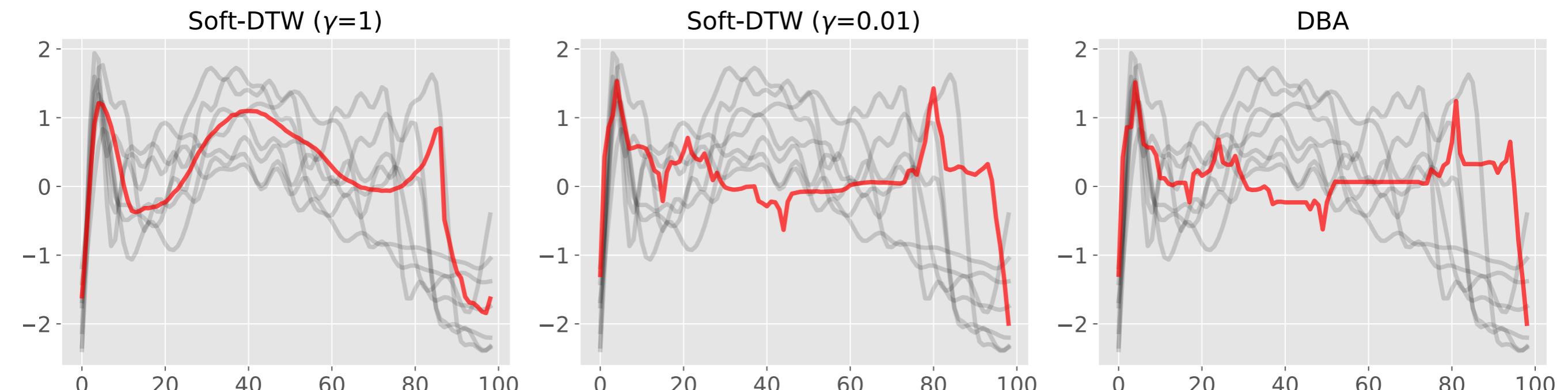
Euclidean loss



Soft-DTW loss ( $\gamma = 1$ )

# sDTW Barycenter

$$\min_{\textcolor{red}{X}} \sum_{j=1} \frac{\lambda_j}{m_j} \text{dtw}_{\gamma}(\textcolor{red}{X}, \textcolor{blue}{Y_j})$$



[DBA] Petitjean et al., A global averaging method for dynamic time warping, with applications to clustering. *Pattern Recognition*, 44 (3):678–693, 2011.

# sDTW Barycenter

$$\min_{\textcolor{red}{X}} \sum_{j=1} \frac{\lambda_j}{m_j} \text{dtw}_{\gamma}(\textcolor{red}{X}, \textcolor{blue}{Y_j})$$

Table 1. Percentage of the datasets on which the proposed soft-DTW barycenter is achieving lower DTW loss (Equation (4) with  $\gamma = 0$ ) than competing methods.

	Random initialization	Euclidean mean initialization
<b>Comparison with DBA</b>		
$\gamma = 1$	40.51%	3.80%
$\gamma = 0.1$	93.67%	46.83%
$\gamma = 0.01$	100%	79.75%
$\gamma = 0.001$	97.47%	89.87%
<b>Comparison with subgradient method</b>		
$\gamma = 1$	96.20%	35.44%
$\gamma = 0.1$	97.47%	72.15%
$\gamma = 0.01$	97.47%	92.41%
$\gamma = 0.001$	97.47%	97.47%

# sDTW Barycenter

$$\min_{\textcolor{red}{X}} \sum_{j=1} \frac{\lambda_j}{m_j} \text{dtw}_{\gamma}(\textcolor{red}{X}, \textcolor{blue}{Y}_j)$$

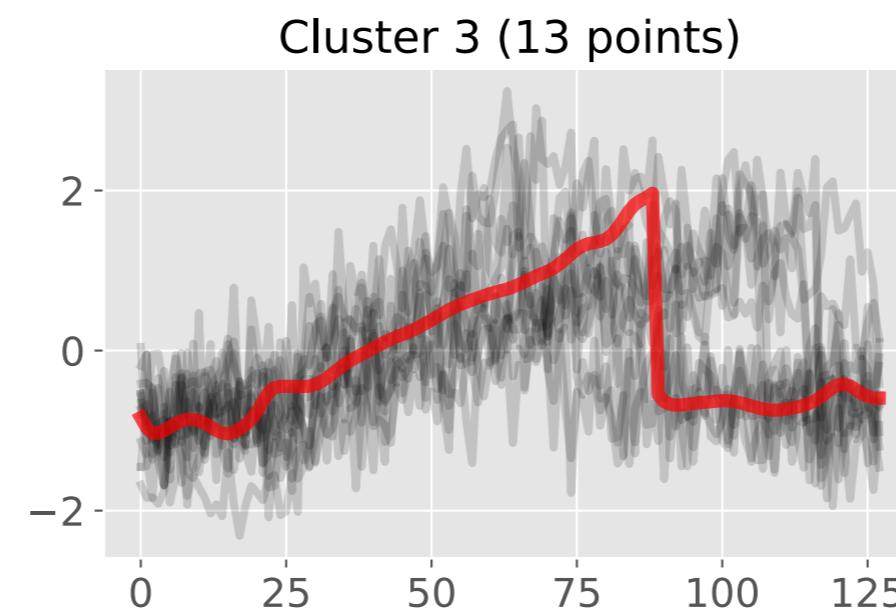
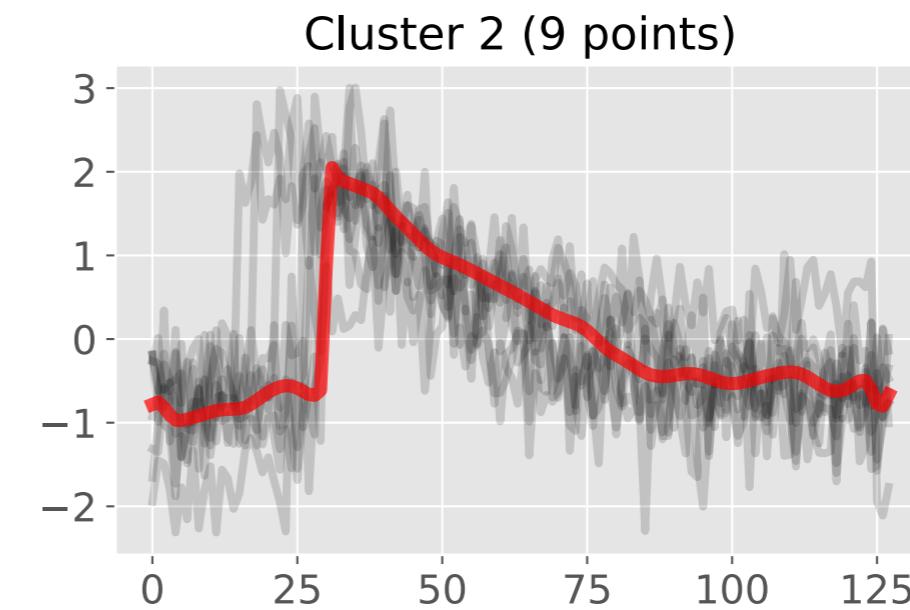
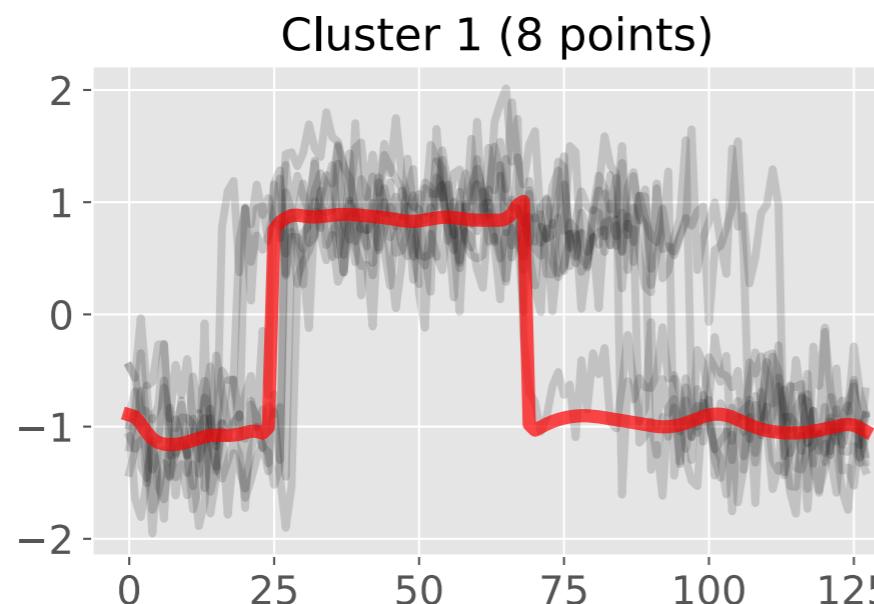
Evaluation performed  
using dtwo

*Table 1.* Percentage of the datasets on which the proposed DTW barycenter is achieving lower DTW loss (Equation 1,  $\gamma = 0$ ) than competing methods.

	Random initialization	Euclidean mean initialization	
<b>Comparison with DBA</b>			
$\gamma = 1$	40.51%	3.80%	% of datasets where soft-dtw is winning
$\gamma = 0.1$	93.67%	46.83%	
$\gamma = 0.01$	100%	79.75%	
$\gamma = 0.001$	97.47%	89.87%	
<b>Comparison with subgradient method</b>			
$\gamma = 1$	96.20%	35.44%	
$\gamma = 0.1$	97.47%	72.15%	
$\gamma = 0.01$	97.47%	92.41%	
$\gamma = 0.001$	97.47%	97.47%	

# sDTW Clustering

$$\min_{\mathbf{X}_1, \dots, \mathbf{X}_k} \sum_{j=1}^N \min_{i=1, \dots, k} \text{dtw}_{\gamma}(\mathbf{X}_i, \mathbf{Y}_j)$$



# sDTW Clustering

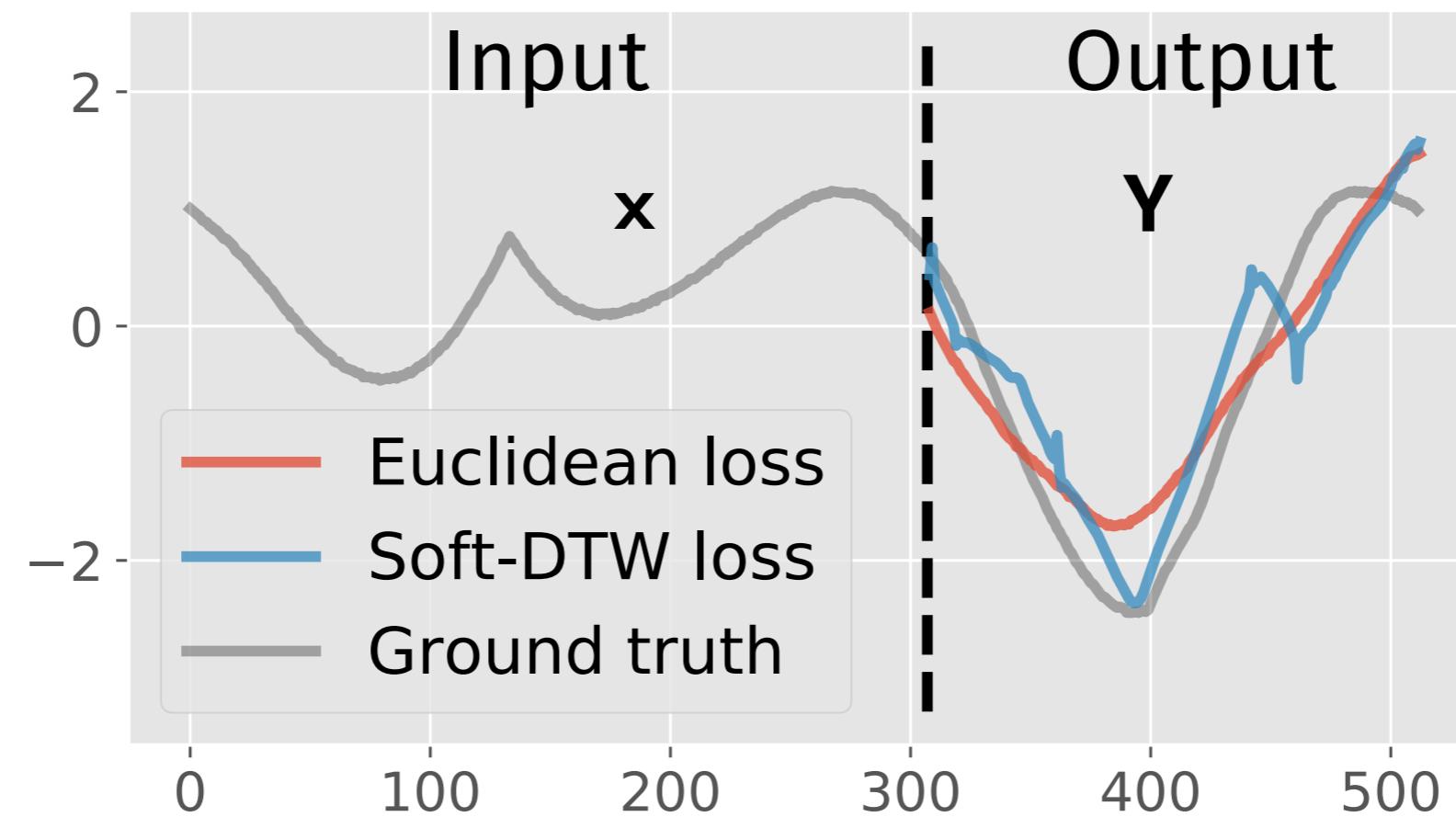
$$\min_{\mathbf{X}_1, \dots, \mathbf{X}_k} \sum_{j=1}^N \min_{i=1, \dots, k} \text{dtw}_{\gamma}(\mathbf{X}_i, \mathbf{Y}_j)$$

Table 2. Percentage of the datasets on which the proposed soft-DTW based ***k*-means** is achieving lower DTW loss (Equation (5) with  $\gamma = 0$ ) than competing methods.

	Random initialization	Euclidean mean initialization
<b>Comparison with DBA</b>		
$\gamma = 1$	15.78%	29.31%
$\gamma = 0.1$	24.56%	24.13%
$\gamma = 0.01$	59.64%	55.17%
$\gamma = 0.001$	77.19%	68.97%
<b>Comparison with subgradient method</b>		
$\gamma = 1$	42.10%	46.44%
$\gamma = 0.1$	57.89%	50%
$\gamma = 0.01$	76.43%	65.52%
$\gamma = 0.001$	96.49%	84.48%

# sDTW Prediction Loss

$$\min_{\theta} \sum_{i=1}^N \frac{1}{m_i} \text{dtw}_{\gamma}(f_{\theta}(x_i), Y_i)$$



# sDTW Prediction Loss

$$\min_{\theta} \sum_{i=1}^N \frac{1}{m_i} \text{dtw}_{\gamma}(f_{\theta}(x_i), Y_i)$$

Table 3. Averaged rank obtained by a multi-layer perceptron (MLP) under Euclidean and soft-DTW losses. Euclidean initialization means that we initialize the MLP trained with soft-DTW loss by the solution of the MLP trained with Euclidean loss.

Training loss	Random initialization	Euclidean initialization	
<b>When evaluating with DTW loss (<math>\text{dtw}_0</math>)</b>			
Euclidean	3.46	4.21	
soft-DTW ( $\gamma = 1$ )	3.55	3.96	
soft-DTW ( $\gamma = 0.1$ )	3.33	3.42	averaged rank
soft-DTW ( $\gamma = 0.01$ )	2.79	2.12	
soft-DTW ( $\gamma = 0.001$ )	<b>1.87</b>	<b>1.29</b>	
<b>When evaluating with Euclidean loss</b>			
Euclidean	<b>1.05</b>	<b>1.70</b>	
soft-DTW ( $\gamma = 1$ )	2.41	2.99	
soft-DTW ( $\gamma = 0.1$ )	3.42	3.38	
soft-DTW ( $\gamma = 0.01$ )	4.13	3.64	
soft-DTW ( $\gamma = 0.001$ )	3.99	3.29	

# Summary

---

- **Dynamic Time Warping** is a natural and flexible discrepancy to compare time series, yet it is **non-differentiable**
- **Soft-DTW** is a differentiable approximation, with better convexity properties
- Using **soft-DTW** typically results in better minima, even when measured with the original DTW
- Python code available on

<https://github.com/mblondel/soft-dtw>