

9 February 2024

Edit Distance

## Plan

\* Edit Distance

\* Announcements

\* ED Dynamic Program

↳ Recurrence

↳ Reconstructing the sequence of edits.

↳ Space Complexity?

# Genomics

\* Genome Sequences  $S \in \{A, C, G, T\}^*$

$S_{\text{cat}} =$  A C C G A T C G A T ...

$S_{\text{dog}} =$  A C G G A A T C G G T ...

# Genomics

\* Genome Sequences  $S \in \{A, C, G, T\}^*$

$S_{\text{cat}} =$  A C C G A T C G A T ...

$S_{\text{dog}} =$  A C G G A A T C G G T ...

How related are two species?

↳ How "similar" are their genomes?

# Genomics

\* Genome Sequences  $S \in \{A, C, G, T\}^*$

$S_{\text{cat}} =$  A C C G A T C G A T ...

$S_{\text{dog}} =$  A C G G A A T C G G T ...

How related are two species?

↳ How "similar" are their genomes?

BLAST: Basic Local Alignment Search Tool > 110 K Citations

# Edit Distance

How easily can we edit  
string  $S$  into string  $t$ ?

# Edit Distance

How easily can we edit string  $S$  into string  $t$ ?

## Allowed Edits

- \* Insertion / Deletion
- \* Changes

~~EDIT~~  
↓ ↓ ↓  
DIST

~~MINIMUM~~  
↓ ↓  
MAXIMUM

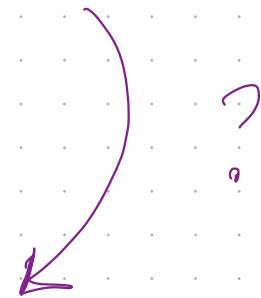
~~SNOW~~  
NO

# Edit Distance

How easily can we edit  
string  $S$  into string  $t$ ?

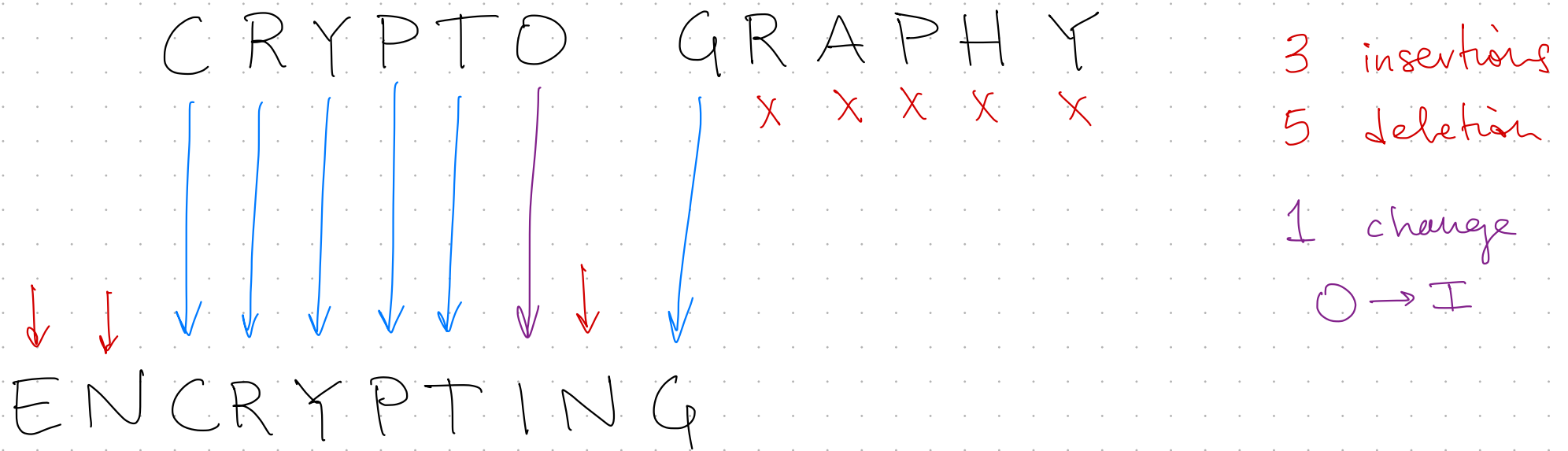
CRYPTOGRAPHY

ENCRYPTING



# Edit Distance

How easily can we edit string  $S$  into string  $t$ ?





Given two strings S and T

Compute minimum cost edits from  $S \rightarrow T$

insertion / deletion :  $\gamma$   
"changing" a to b :  $\Delta_{ab}$

S A G G C T A A T C . . .

T G A G C T A A G C C . . .

Changing a to a  
costs  $\Delta_{aa} = 0$ .

# Announcements

\* HW2 due last night. Late deadline Sunday.

\* HW3 out after lecture.

\* Prelim Tues Feb 20 7:30 pm

↳ Last call to schedule a make-up

See Ed post from Sara Perkins

Given two strings  $S$  and  $T$

Compute minimum cost edits from  $S \rightarrow T$

insertion / deletion :  $\gamma$

changing  $a$  to  $b$  :  $\Delta_{ab}$

$S$  AGGCTAATC . . .

$T$  GAGCTAAGCC . . .

Given two strings  $S$  and  $T$

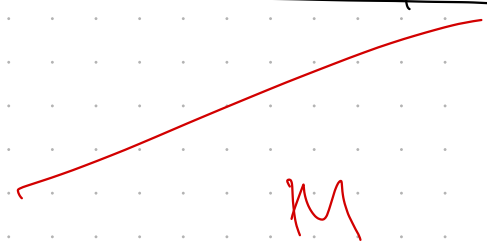
Compute minimum cost edits from  $S \rightarrow T$

insertion / deletion :  $\gamma$

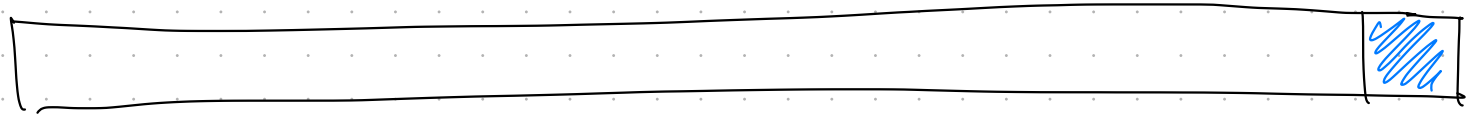
changing  $a$  to  $b$  :  $\Delta_{ab}$

$S$  [AGGCTAATC . . . ]   $n$

$T$  [GAGCTAAGCC . . . ]   $m$



S



$S_n$

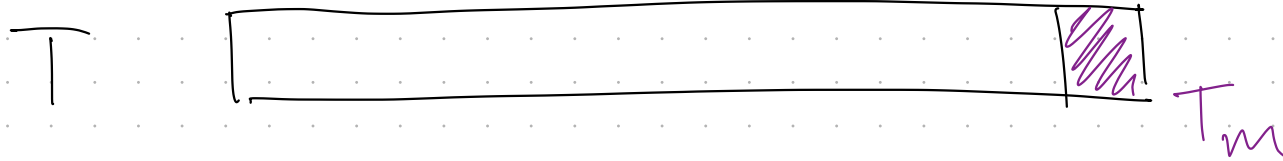
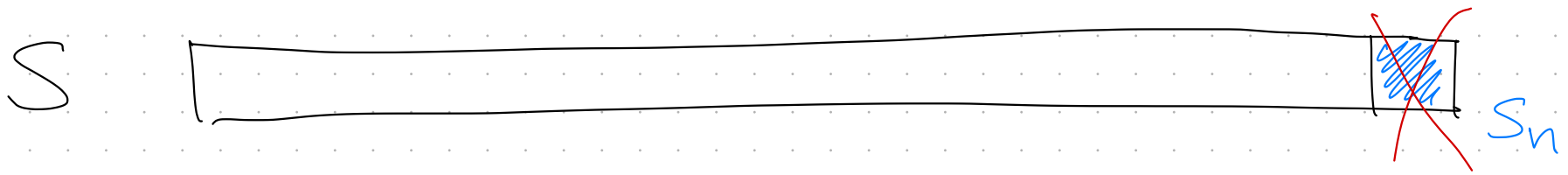
T



$T_m$

Fact, Fix an optimal set of edits from S to T.

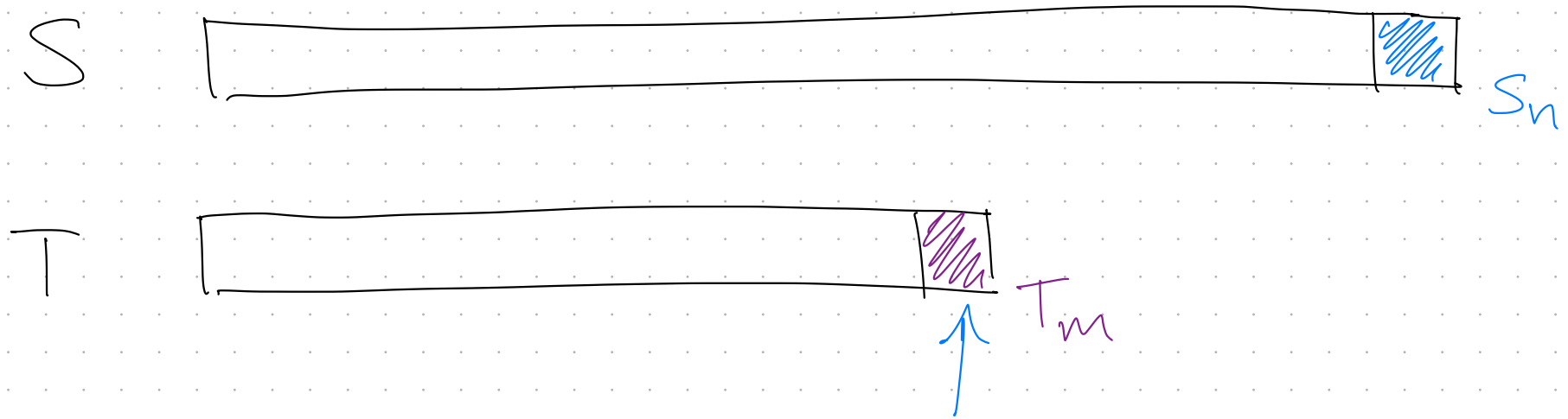
One of the following is true:



Fact, Fix an optimal set of edits from S to T.

One of the following is true:

\*  $S_n$  is ~~deleted~~ at cost  $\lambda$

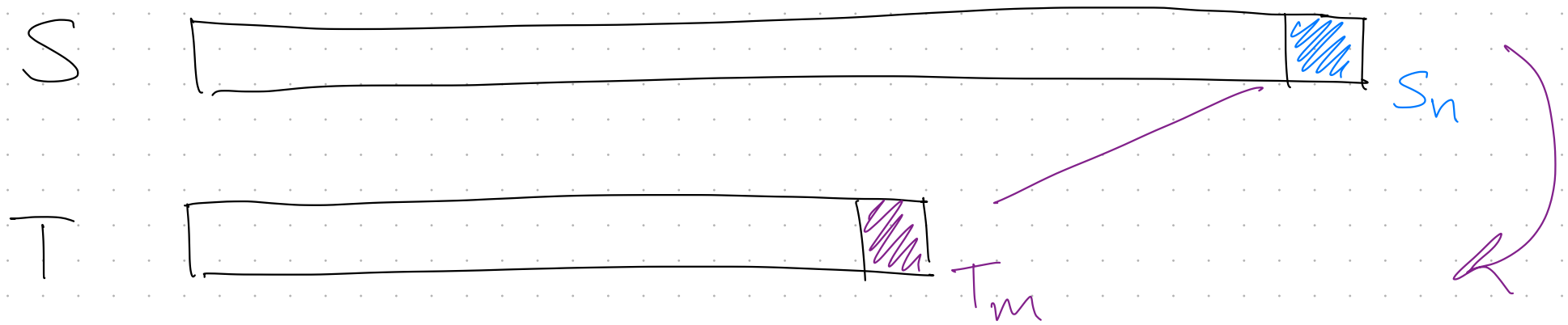


Fact, Fix an optimal set of edits from S to T.

One of the following is true:

\*  $S_n$  is ~~deleted~~ at cost  $\gamma$

\*  $T_m$  is ~~inserted~~ at cost  $\gamma$



Fact, Fix an optimal set of edits from S to T.

One of the following is true:

\*  $S_n$  is ~~deleted~~ at cost  $\gamma$

\*  $T_m$  is inserted at cost  $\gamma$

\*  $S_n$  "changes" to  $T_m$  at cost  $\Delta_{S_n T_m}$



# The Edit Distance Recurrence




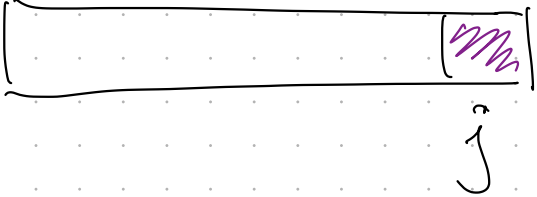
$ED(i, j)$   
↓  
Edit distance  
between prefixes  
 $S[1:i]$  and  $T[1:j]$

## Possibilities

- \*  $S_i$  changes to  $T_j$
- \*  $S_i$  deleted
- \*  $T_j$  inserted

# The Edit Distance Recurrence

$S[1:i]$  

$T[1:j]$  

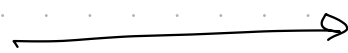
$ED(i, j)$   
minimizes

## Possibilities

\*  $S_i$  changes to  $T_j$

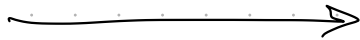
$$ED(i-1, j-1) + \Delta_{S_i T_j}$$

\*  $S_i$  deleted



$$ED(i-1, j) + \gamma$$

\*  $T_j$  inserted



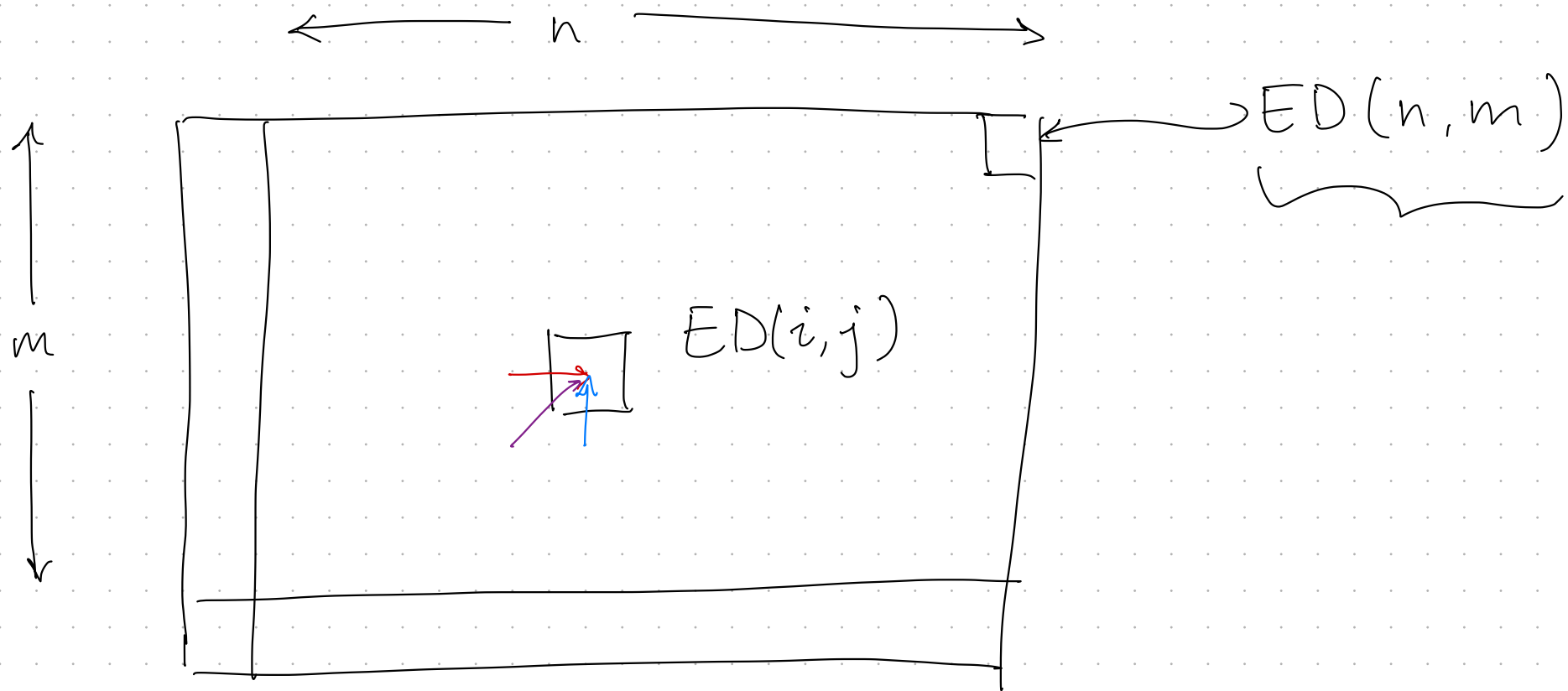
$$ED(i, j-1) + \gamma$$

$$ED(i, j) = \min \left\{ \begin{array}{l} ED(i-1, j-1) + \Delta s_i T_j \\ ED(i-1, j) + \gamma \\ ED(i, j-1) + \gamma \end{array} \right.$$

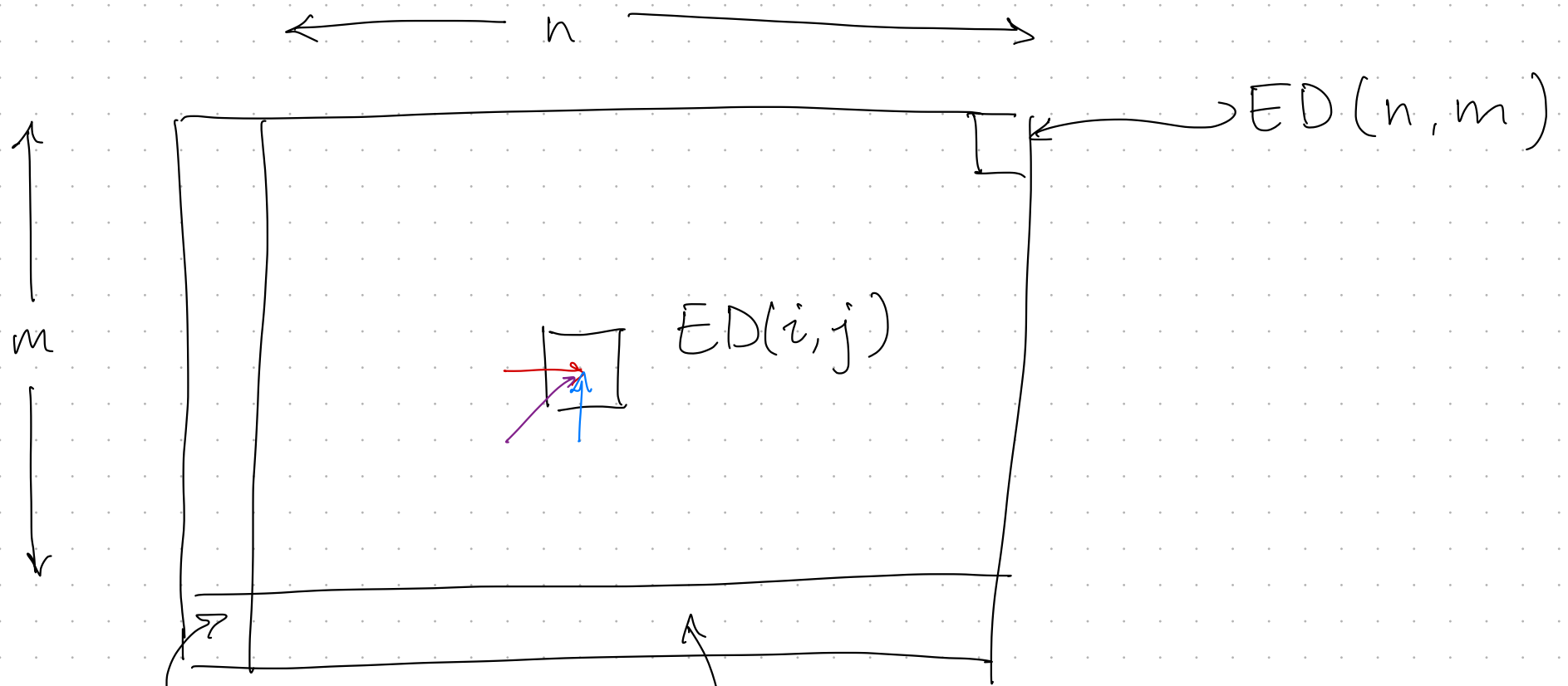
---

Dynamic Programming Table?

Base Cases?



$$ED(i, j) = \min \left\{ \begin{array}{l} ED(i-1, j-1) + \Delta_{s_i t_j} \\ ED(i-1, j) + \gamma \\ ED(i, j-1) + \gamma \end{array} \right.$$



$$ED(0, 0) = 0$$

$$ED(i, 0) = i \cdot \gamma$$

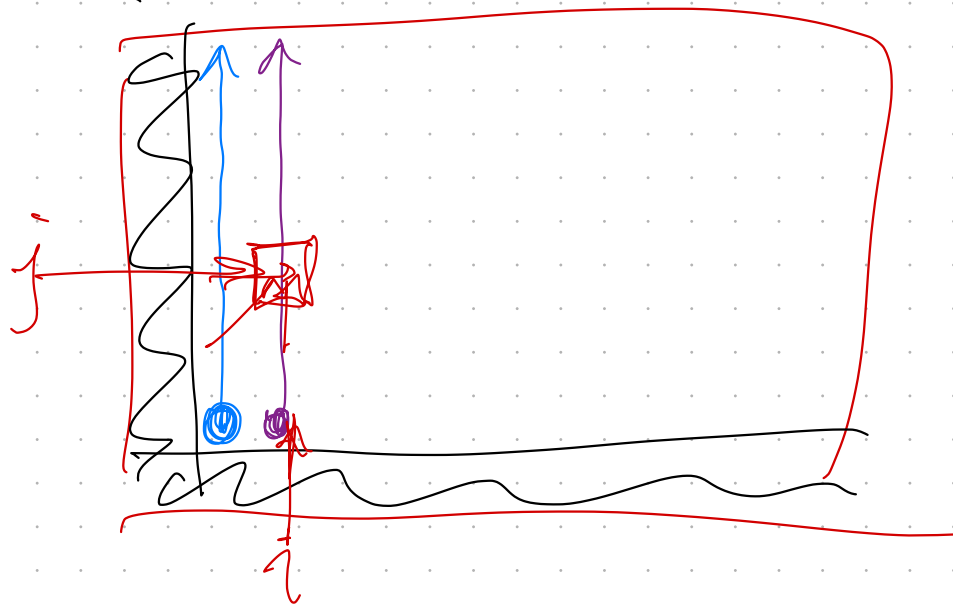
$$ED(0, j) = j \cdot \delta$$

# Edit Distance Algorithm

$$ED(0,0) = 0$$

$$ED(i,0) = i \cdot \gamma$$

$$ED(0,j) = j \cdot \gamma$$



For  $i = 1 \rightarrow n$

For  $j = 1 \rightarrow m$

$$ED(i,j) = \min \left\{ \begin{array}{l} ED(i-1, j-1) \\ \quad + \Delta_{s_i T_j} \end{array} , \begin{array}{l} ED(i-1, j) \\ \quad + \gamma \end{array} , \begin{array}{l} ED(i, j-1) \\ \quad + \gamma \end{array} \right\}$$

return  $ED(n, m)$

# Edit Distance Algorithm

$$ED(0,0) = 0$$

$$ED(i,0) = i \cdot \gamma$$

$$ED(0,j) = j \cdot \gamma$$

For  $i = 1 \rightarrow n$

For  $j = 1 \rightarrow m$

$$ED(i,j) = \min \left\{ \begin{array}{l} ED(i-1, j-1) \\ + \Delta_{s_i T_j} \end{array} , \begin{array}{l} ED(i-1, j) \\ + \gamma \end{array} , \begin{array}{l} ED(i, j-1) \\ + \gamma \end{array} \right\}$$

return  $ED(n, m)$

Running Time

-  $1+n+m$  initialization

-  $n$  outer iterations

$m$  inner iterations

$O(1)$  per iteration

# Edit Distance Algorithm

$$ED(0,0) = 0$$

$$ED(i,0) = i \cdot \gamma$$

$$ED(0,j) = j \cdot \gamma$$

For  $i = 1 \rightarrow n$

For  $j = 1 \rightarrow m$

$$ED(i,j) = \min \left\{ \begin{array}{l} ED(i-1,j-1) \\ + \Delta_{s_i T_j} \end{array} , \begin{array}{l} ED(i-1,j) \\ + \gamma \end{array} , \begin{array}{l} ED(i,j-1) \\ + \gamma \end{array} \right\}$$

return  $ED(n,m)$

Running Time

-  $1+n+m$  initialization

-  $n$  outer iterations

$m$  inner iterations

$O(1)$  per iteration

Space

DP Table has  $(m+1) \times (n+1)$  cells

$O(mn)$



Edit Distance  $\equiv$  Sequence Alignment

- - CRYPT-OGRAPHY  
| | | | | | |  
ENCRYPTING - - - -

An alignment is a non-crossing matching.

Edit Distance  $\equiv$  Sequence Alignment

insertions



deletions

An alignment is a non-crossing matching.

Edit Distance  $\equiv$  Sequence Alignment

- - CRYPT-OGRAPHY  
| | | | | |  
ENCRYPTING - - - -  
changes

An alignment is a non-crossing matching.

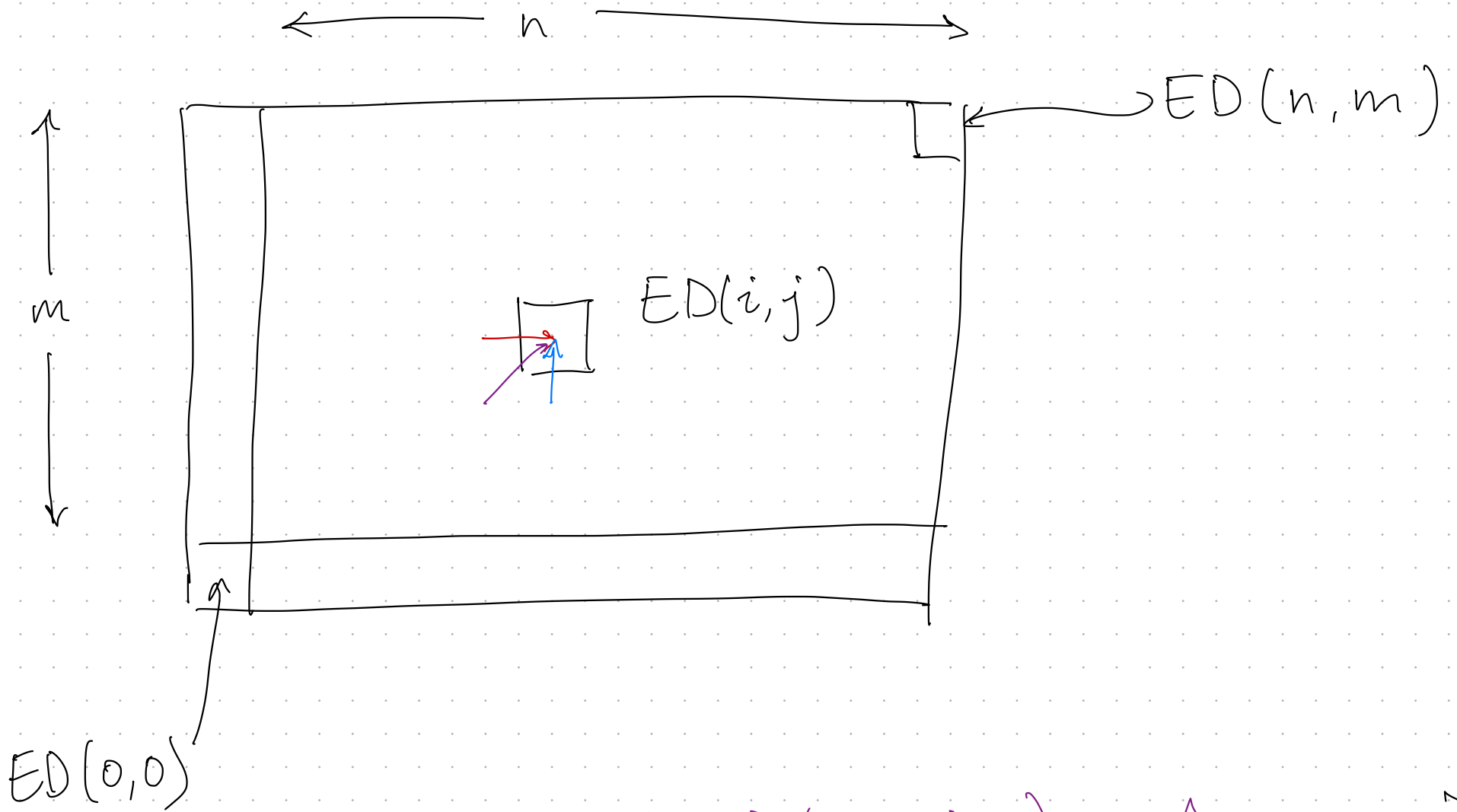
Edit Distance  $\equiv$  Sequence Alignment

CRYPTOGRAPHY  
| | | | | | |  
ENCRYPTING

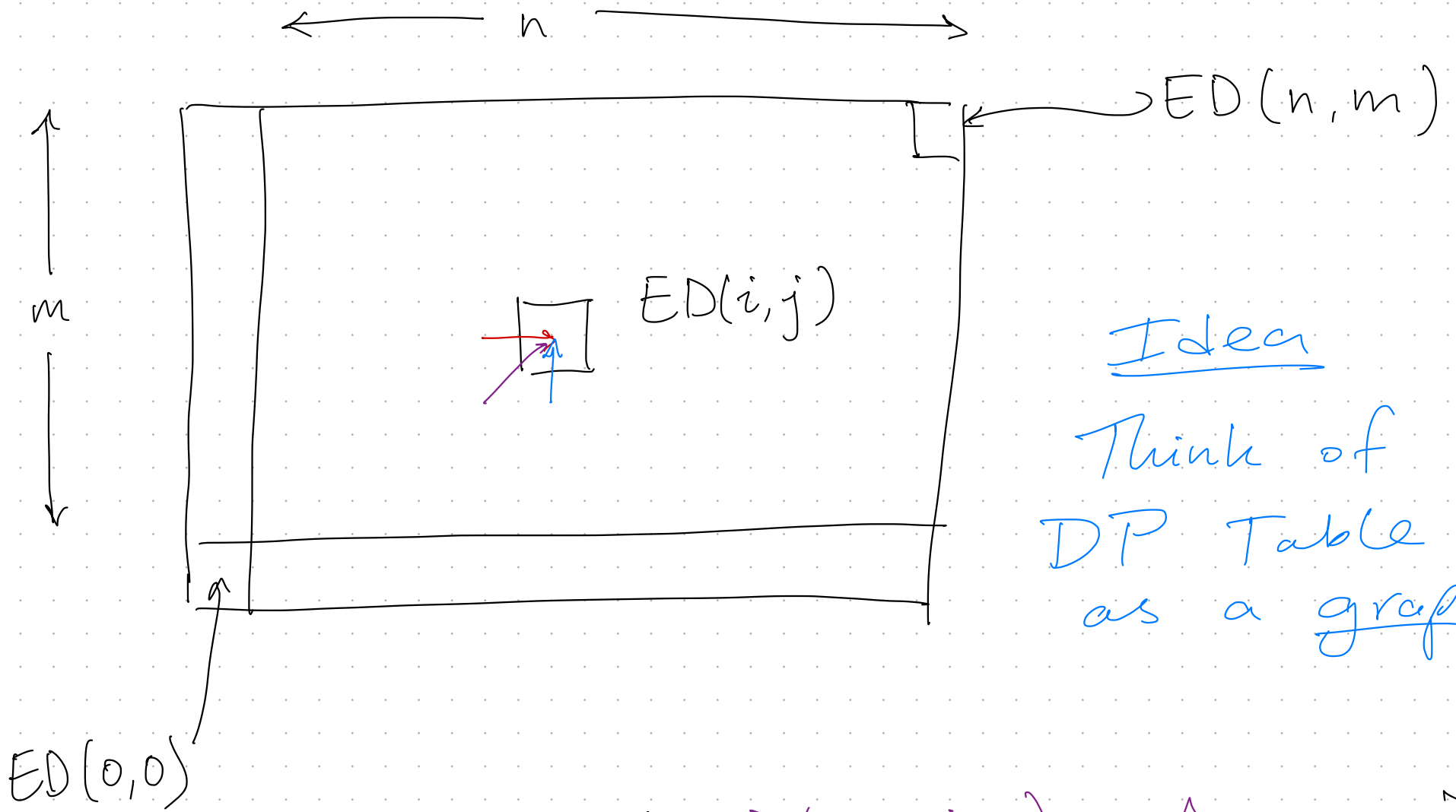
An alignment is a non-crossing matching.

ED Algorithm computed the distance

Can we compute the alignment?

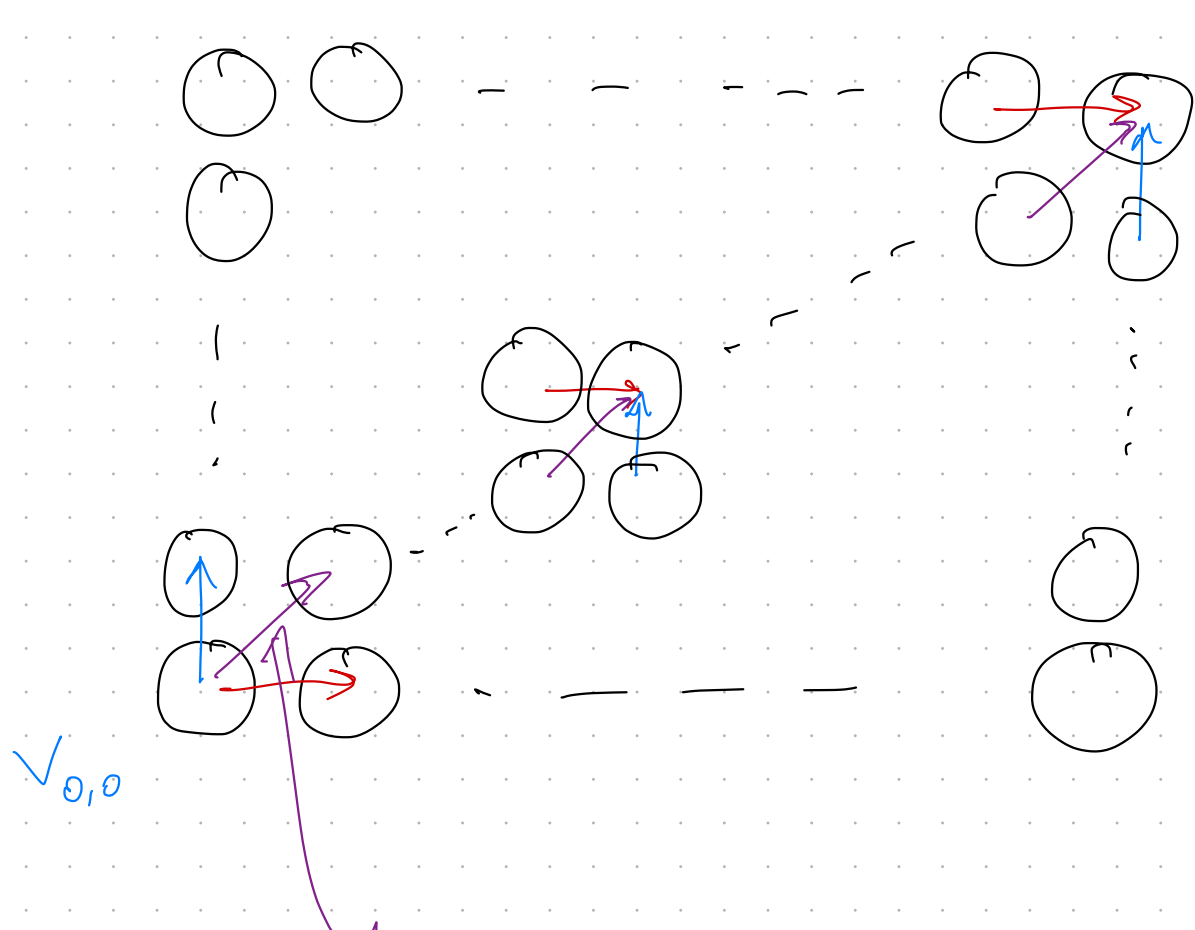


$$ED(i, j) = \min \left\{ \begin{array}{l} ED(i-1, j-1) + \Delta_{s_i t_j} \\ ED(i-1, j) + \gamma \\ ED(i, j-1) + \gamma \end{array} \right.$$



Idea  
Think of  
DP Table  
as a graph

$$ED(i,j) = \min \left\{ \begin{array}{l} ED(i-1, j-1) + \Delta_{s_i T_j} \\ ED(i-1, j) + \gamma \\ ED(i, j-1) + \gamma \end{array} \right.$$



$V_{n,m}$

Vertices  
 $ED(i,j) \equiv v_{i,j}$

Edges  
 correspond to  
 cost of edit

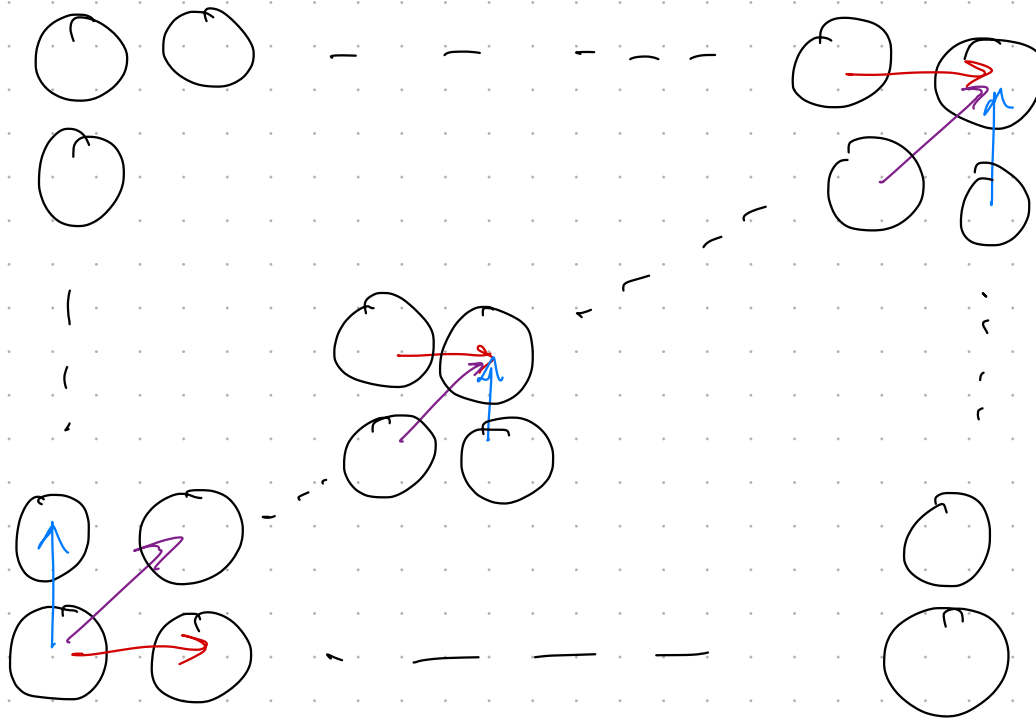
$V_{0,0}$

$\Delta_{s_i t_j}$

$ED(i,j) = \min$

$ED(i-1, j-1) + \Delta_{s_i t_j}$  ,  
 $ED(i-1, j) + \gamma$  ,  
 $ED(i, j-1) + \gamma$

$$S = V_{0,0}$$



$$V_{n,m} = t$$

$$\frac{\text{Vertices}}{ED(i,j)} \equiv v_{i,j}$$

Edges  
correspond to  
cost of edit

---

Minimum Cost Alignment

≡

Shortest st-Path in ED  
Graph



Theorem. Shortest path from  $v_{00}$  to  $v_{ij}$   
equals  $ED(i, j)$ .

---

Pf. By induction on  $i+j$ .

Base Case. Shortest path from  $v_{00} \rightarrow v_{00} = 0$   
 $= ED(0, 0)$

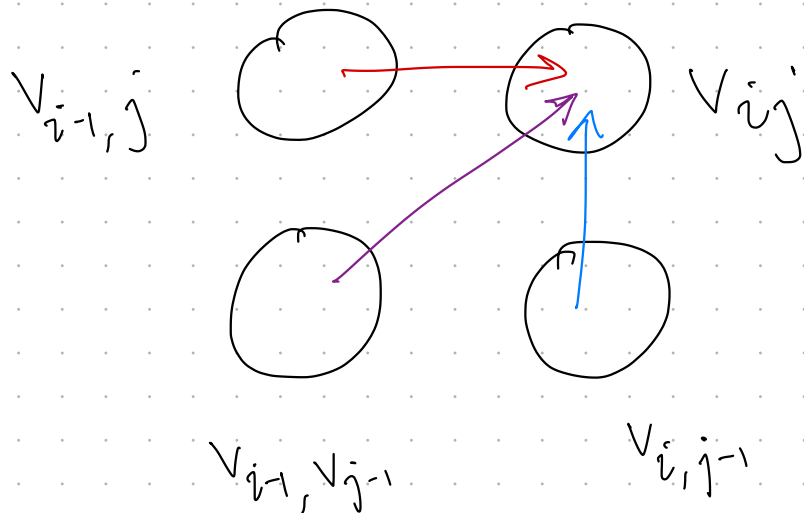
Theorem. Shortest path from  $v_{00}$  to  $v_{ij}$   
equals  $ED(i, j)$ .

---

Pf. By induction on  $i+j$ .

Base Case. Shortest path from  $v_{00} \rightarrow v_{00} = 0$   
 $= ED(0, 0)$

Inductive Step. Suppose shortest path to  $v_{kl} = ED(k, l)$   
for all  $k+l < i+j$ .



Theorem. Shortest path from  $v_{00}$  to  $v_{ij}$  equals  $ED(i, j)$ .

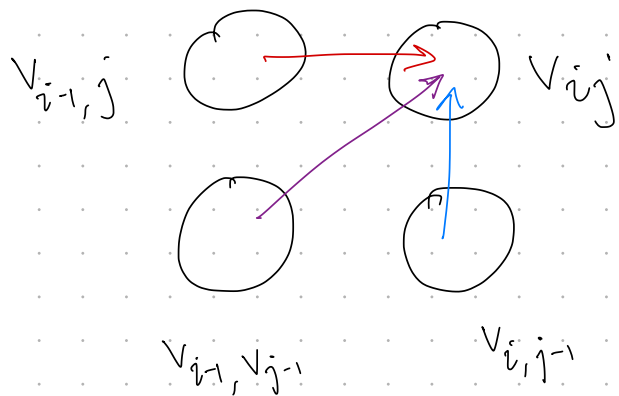
---

Pf. By induction on  $i+j$ .

Base Case. Shortest path from  $v_{00} \rightarrow v_{00} = 0$   
 $= ED(0, 0)$

Inductive Step. Suppose shortest path to  $v_{kl} = ED(k, l)$  for all  $k+l < i+j$ .

$$SP(v_{00} \rightarrow v_{ij}) = \min \left\{ \begin{array}{l} SP(v_{00} \rightarrow v_{i-1, j-1}) + \Delta_{S_i T_j} \\ SP(v_{00} \rightarrow v_{i-1, j}) + \gamma \\ SP(v_{00} \rightarrow v_{i, j-1}) + \gamma \end{array} \right\}$$



Theorem. Shortest path from  $v_{00}$  to  $v_{ij}$   
equals  $ED(i, j)$ .

---

Pf. By induction on  $i+j$ .

Base Case. Shortest path from  $v_{00} \rightarrow v_{00} = 0$   
 $= ED(0, 0)$

Inductive Step. Suppose shortest path to  $v_{kl} = ED(k, l)$   
for all  $k+l < i+j$ .

$$SP(v_{00} \rightarrow v_{ij}) = \min \left\{ \begin{array}{l} ED(i-1, j-1) + \Delta_{s_{i-1}t_j} \\ ED(i-1, j) + \gamma \\ ED(i, j-1) + \gamma \end{array} \right\}$$

Theorem. Shortest path from  $v_{00}$  to  $v_{ij}$   
equals  $ED(i, j)$ .

---

Pf. By induction on  $i+j$ .

Base Case. Shortest path from  $v_{00} \rightarrow v_{00} = 0$   
 $= ED(0, 0)$

Inductive Step. Suppose shortest path to  $v_{kl} = ED(k, l)$   
for all  $k+l < i+j$ .

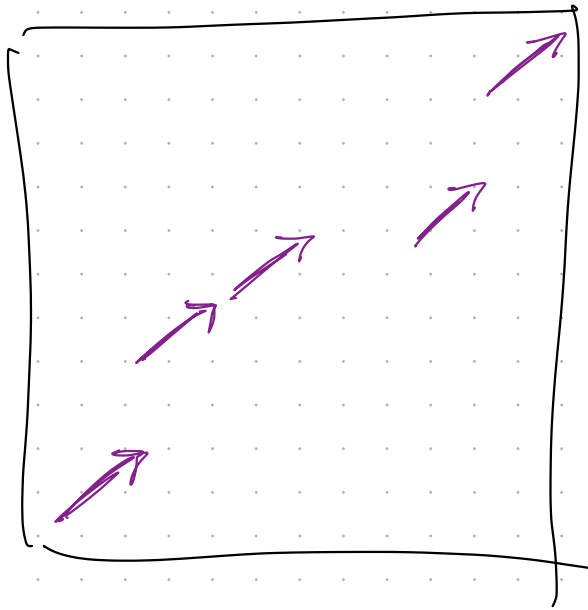
$$SP(v_{00} \rightarrow v_{ij}) = \min \left\{ \begin{array}{l} ED(i-1, j-1) + \Delta_{s_{i-1}t_j} \\ ED(i-1, j) + \gamma \\ ED(i, j-1) + \gamma \end{array} \right\}$$

$$= ED(i, j)$$

Theorem. Shortest path from  $v_{00}$  to  $v_{ij}$   
equals  $ED(i, j)$ .

---

Corollary. Diagonal edges on shortest  
path from  $v_{00} \rightarrow v_{nm}$  correspond to alignment.



A	C	G	T	T	C	A
↓	↘	↘	↗	↓		
A	A	C	C	C	G	A

# Edit Distance Algorithm

$$ED(0,0) = 0$$

$$ED(i,0) = i \cdot \gamma$$

$$ED(0,j) = j \cdot \gamma$$

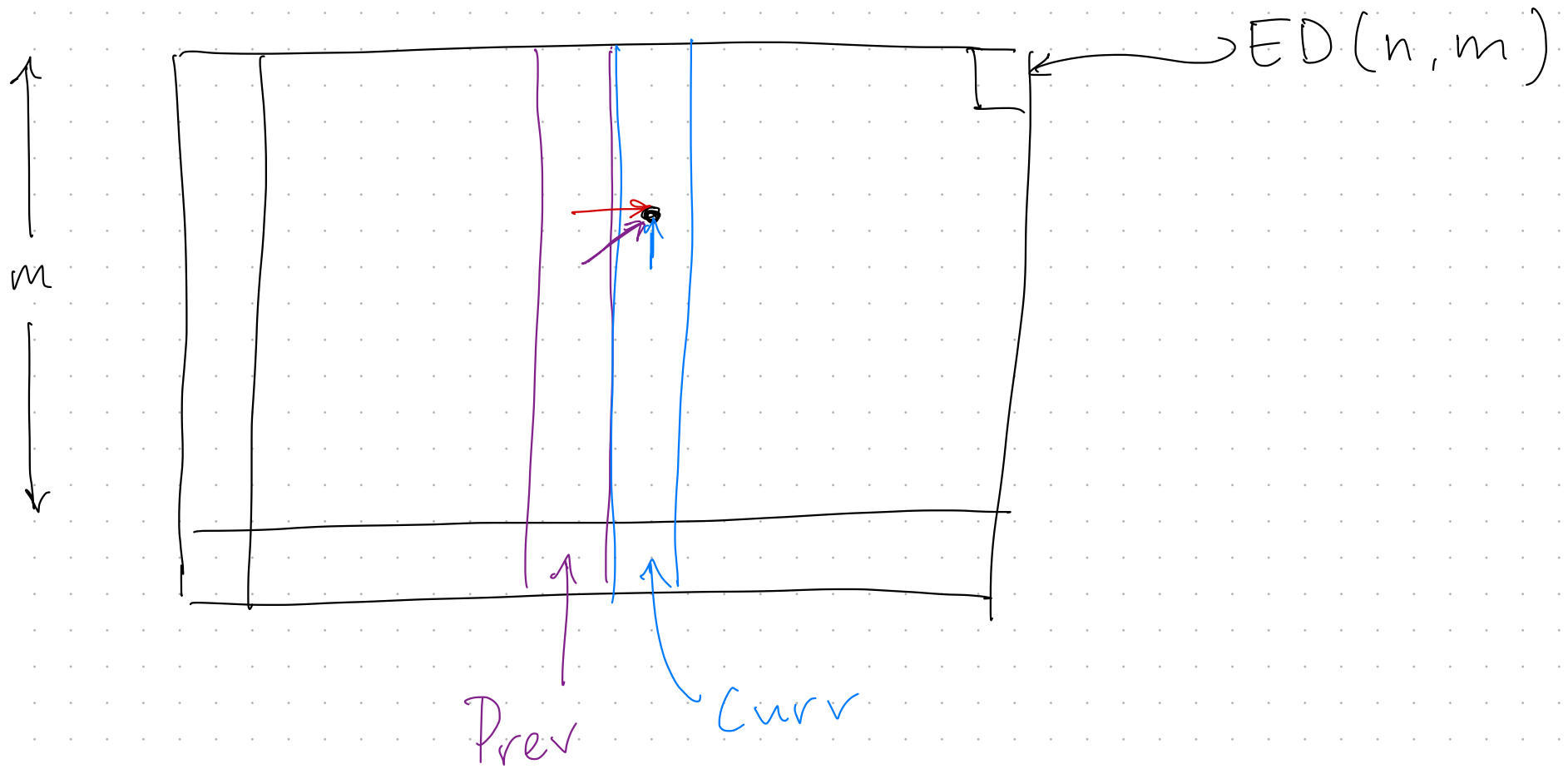
How much space did we really need?

For  $i = 1 \rightarrow n$

For  $j = 1 \rightarrow m$

$$ED(i,j) = \min \left\{ \begin{array}{l} ED(i-1, j-1) \\ + \Delta_{s_i T_j} \end{array} , \begin{array}{l} ED(i-1, j) \\ + \gamma \end{array} , \begin{array}{l} ED(i, j-1) \\ + \gamma \end{array} \right\}$$

return  $ED(n, m)$



$$Curr(j) = \min \left\{ \begin{array}{l} Prev(j-1) + \Delta_{s_i T_j} \\ Prev(j) + \gamma \\ Curr(j-1) + \gamma \end{array} \right\}$$



# Linear Space ED.

$$\text{Prev}(j) = j \cdot \gamma$$

// m-entry 1D arrays

$$\text{Curr}(j) = 0$$

For  $i = 1 \rightarrow n$

$$\text{Curr}(0) = i \cdot \gamma$$

For  $j = 1 \rightarrow m$

$$\text{Curr}(j) = \min \left\{ \begin{array}{l} \text{Prev}(j-1) \\ + \Delta_{s_i T_j} \end{array} , \begin{array}{l} \text{Prev}(j) \\ + \gamma \end{array} , \begin{array}{l} \text{Curr}(j-1) \\ + \gamma \end{array} \right\}$$

Prev  $\leftarrow$  Curr.

return Curr(m).