

String matching: definition of the problem (text,pattern)

- **Exact matching:** depends on what we have: text or patterns
 - **The patterns** ---> Data structures for the patterns
 - 1 pattern ---> The algorithm depends on $|p|$ and $|\Sigma|$
 - k patterns ---> The algorithm depends on k, $|p|$ and $|\Sigma|$
 - Extensions
 - Regular Expressions
 - **The text** ----> Data structure for the text (suffix tree, ...)
- **Approximate matching:**
 - Dynamic programming
 - Sequence alignment (pairwise and multiple)
 - Sequence assembly: hash algorithm
- **Probabilistic search:** Hidden Markov Models

Approximate string matching

For instance, given the sequence

CTACTACTACGTGACTAATACTGATCGTAGCTAC...

search for the pattern ACTGA allowing one error...

... but what is the meaning of “one error”?

We accept three types of errors:

- 1. Mismatch: ACCG**T**GAT ACCG**A**GAT
 - 2. Insertion: ACCGTGAT ACCG**A**TGAT
 - 3. Deletion: ACCG**T**GAT ACCGGAT
- } Indel

The **edit distance** d between two strings is the minimum number of substitutions, insertions and deletions needed to transform the first string into the second one

$$d(\text{ACT}, \text{ACT}) =$$

$$d(\text{ACT},) =$$

$$d(\text{ACT}, \text{AC}) =$$

$$d(\text{AC}, \text{ATC}) =$$

$$d(\text{ACT}, \text{C}) =$$

$$d(\text{ACTTG}, \text{ATCTG}) =$$

We accept three types of errors:

1. Mismatch: ACCG**T**GAT ACCG**A**GAT
 2. Insertion: ACCG**T**GAT ACCG**A**TGAT
 3. Deletion: ACCG**T**GAT ACCG**G**GAT
- } Indel

The **edit distance** d between two strings is the minimum number of substitutions, insertions and deletions needed to transform the first string into the second one

$$d(\text{ACT}, \text{ACT}) = 0$$

$$d(\text{ACT},) = 3$$

$$d(\text{ACT}, \text{AC}) = 1$$

$$d(\text{AC}, \text{ATC}) = 1$$

$$d(\text{ACT}, \text{C}) = 2$$

$$d(\text{ACTTG}, \text{ATCTG}) = 2$$

The Edit distance is related with the best alignment of strings

Given

$$d(\text{ACT}, \text{ACT})=0 \quad d(\text{ACT}, \text{AC})=1 \quad d(\text{ACTTG}, \text{ATCTG})=2$$

which is the best alignment in every case?

- ACT and ACT :
ACT
ACT
- ACT and AT:
ACT
A - T
- ACTTG and ATCTG:
ACTTG ACT - TG
ATCTG A - TCTG

Then, the alignment suggest the substitutions, insertions and deletions to transform one string into the other

But which is the distance between the strings

ACGCTATGCTATACG and ACGGTAGTGACGC?

... and the best alignment between them?

1966 was the first time this problem was discussed...

and the algorithm was proposed in 1968, 1970, ...

using the technique called “**Dynamic programming**”

Edit distance and alignment of strings

| | C | T | A | C | T | A | C | G | T |
|---|---|---|---|---|---|---|---|---|---|
| A | | | | | | | | | |
| C | | | | | | | | | |
| T | | | | | | | | | |
| G | | | | | | | | | |
| A | | | | | | | | | |

A grid for dynamic programming. The top row contains the string 'CTACTACTACGT' in teal. The left column contains the string 'ACTGA' in teal. A red square is located at the intersection of the 'C' row and the 5th column (the second 'T' in the top string).

Edit distance and alignment of strings

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|--|
| | C | T | A | C | T | A | C | G | T | |
| A | | | | | | | | | | |
| C | | | | | | | | | | |
| T | | | | | | | | | | |
| G | | | | | | | | | | |
| A | | | | | | | | | | |

The cell contains the distance between AC and CTACT.

Edit distance and alignment of strings

| | | C | T | A | C | T | A | C | G | T | |
|---|---|---|---|---|---|---|---|---|---|---|--|
| | 0 | 1 | ? | | | | | | | | |
| A | | | | | | | | | | | |
| C | | | | | | | | | | | |
| T | | | | | | | | | | | |
| G | | | | | | | | | | | |
| A | | | | | | | | | | | |

-
C

Edit distance and alignment of strings

| | | C | T | A | C | T | A | C | G | T | |
|---|---|---|---|---|---|---|---|---|---|---|--|
| | 0 | 1 | 2 | ? | | | | | | | |
| A | | | | | | | | | | | |
| C | | | | | | | | | | | |
| T | | | | | | | | | | | |
| G | | | | | | | | | | | |
| A | | | | | | | | | | | |

--
CT

Edit distance and alignment of strings

| | | C | T | A | C | T | A | C | G | T | |
|---|---|---|---|---|---|---|---|---|---|-----|--|
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | ... | |
| A | | | | | | | | | | | |
| C | | | | | | | | | | | |
| T | | | | | | | | | | | |
| G | | | | | | | | | | | |
| A | | | | | | | | | | | |

CTACTA

Edit distance and alignment of strings

| | | C | T | A | C | T | A | C | G | T | | |
|---|-----|---|---|---|---|---|---|---|---|-----|--|--|
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | ... | | |
| A | 1 | | | | | | | | | | | |
| C | 2 | | | | | | | | | | | |
| T | 3 | | | | | | | | | | | |
| G | ... | | | | | | | | | | | |
| A | | | | | | | | | | | | |

ACT

Edit distance and alignment of strings

| | | C | T | A | C | T | A | C | G | T |
|---|---|---|---|---|---|---|---|---|---|-----|
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | ... |
| A | 1 | | | | | | | | | |
| C | 2 | | | | | | | | | |
| T | 3 | | | | | | | | | |
| G | | | | | | | | | | |
| A | | | | | | | | | | |

$$\begin{aligned}
 & \text{BA(AC,CTAC)} = \text{best} \\
 & d(\text{AC,CTAC}) = \text{min}
 \end{aligned}
 \left\{ \begin{array}{l}
 \text{BA(AC,CTA)} \begin{array}{l} - \\ C \end{array} & d(\text{AC,CTA})+1 \\
 \text{BA(A,CTA)} \begin{array}{l} C \\ C \end{array} & d(\text{A,CTA}) \\
 \text{BA(A,CTAC)} \begin{array}{l} C \\ - \end{array} & d(\text{A,CTAC})+1
 \end{array} \right.$$

Edit distance and alignment of strings

Connect to

<http://alggen.lsi.upc.es/docencia/ember/leed/Tfc1.htm>

and use the global method.

How this algorithm can be applied

to the approximate search?

to the K-approximate string searching?

K-approximate string searching

| | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|-----|--|
| | C | T | A | C | T | A | C | T | A | C | G | T | A | C | T | G | G | T | G | A | ... | |
| A | | | | | | | | | | | | | | | | | | | | | | |
| C | | | | | | | | | | | | | | | | | | | | | | |
| T | | | | | | | | | | | | | | | | | | | | | | |
| G | | | | | | | | | | | | | | | | | | | | | | |
| A | | | | | | | | | | | | | | | | | | | | | | |

This cell ...

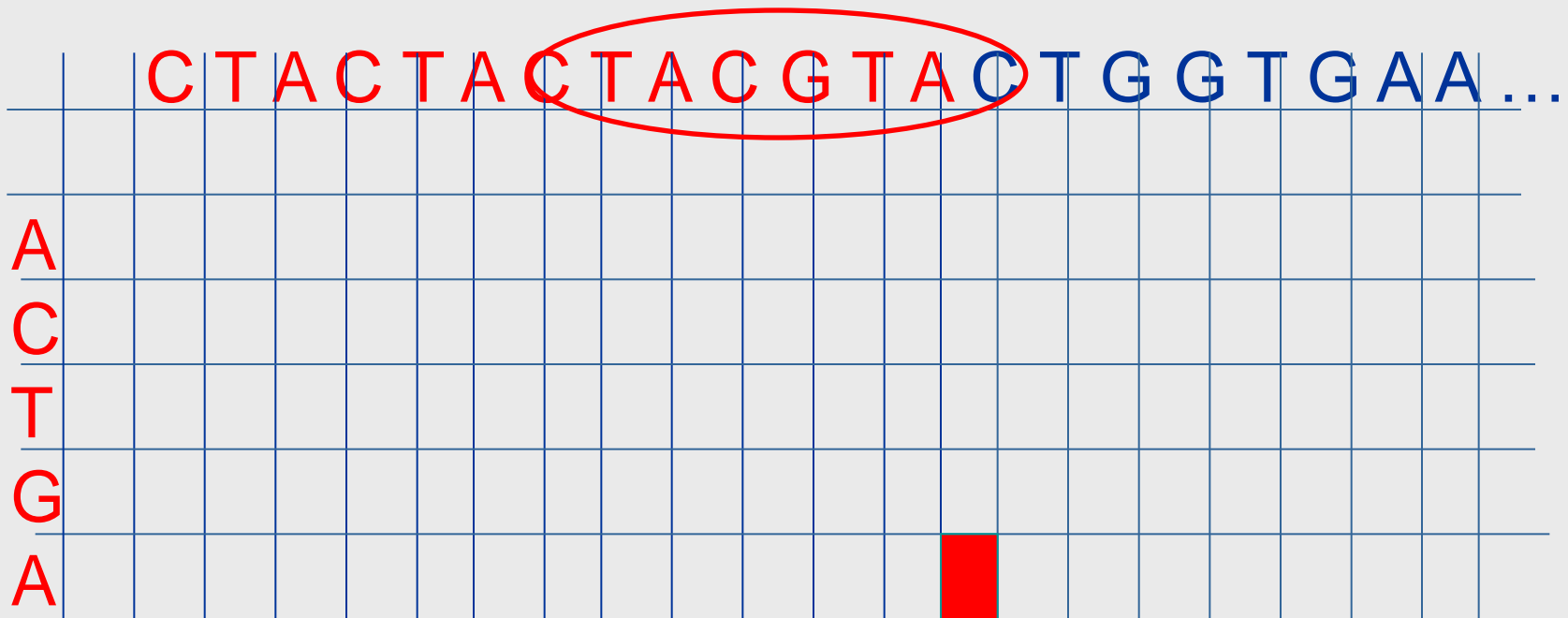
K-approximate string searching

| | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|-----|--|
| | C | T | A | C | T | A | C | T | A | C | G | T | A | C | T | G | G | T | G | A | A | ... | |
| A | | | | | | | | | | | | | | | | | | | | | | | |
| C | | | | | | | | | | | | | | | | | | | | | | | |
| T | | | | | | | | | | | | | | | | | | | | | | | |
| G | | | | | | | | | | | | | | | | | | | | | | | |
| A | | | | | | | | | | | | | | | | | | | | | | | |

This cell gives the distance between (ACTGA, CT...GTA)...

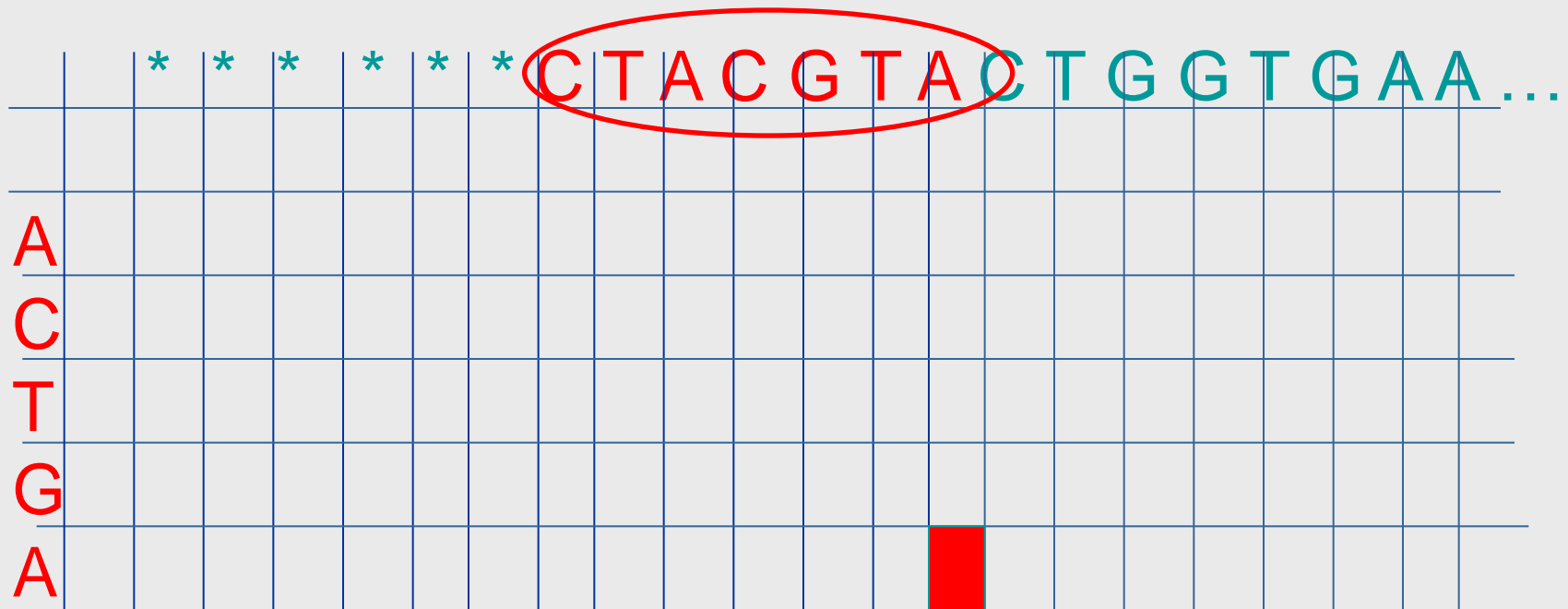
...but we only are interested in the last characters

K-approximate string searching



This cell gives the distance between (ACTGA, CT...GTA)...
...but we only are interested in the last characters

K-approximate string searching



This cell gives the distance between (ACTGA, CT...GTA)...

...but we only are interested in the last characters...

...no matter where they appears in the text, then...

K-approximate string searching



This cell gives the distance between (ACTGA, CT...GTA)...

...but we only are interested in the last characters...

...no matter where they appears in the text, then...

K-approximate string searching



This cell gives the distance between (ACTGA, CT...GTA)...

...but we only are interested in the last characters...

...no matter where they appears in the text, then...

K-approximate string searching

| | | | | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|-----|---|
| | | C | T | A | C | T | A | C | T | A | C | G | T | A | C | T | G | G | T | G | A | A | ... | |
| | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| A | | | | | | | | | | | | | | | | | | | | | | | | |
| C | | | | | | | | | | | | | | | | | | | | | | | | |
| T | | | | | | | | | | | | | | | | | | | | | | | | |
| G | | | | | | | | | | | | | | | | | | | | | | | | |
| A | | | | | | | | | | | | | | | | | | | | | | | | |

This cell gives the distance between (ACTGA, CT...GTA)...

...but we only are interested in the last characters...

...no matter where they appears in the text, then

K-approximate string searching

Connect to

<http://alggen.lsi.upc.es/docencia/ember/leed/Tfc1.htm>

and use the semi-global method.