# Sequence algorithms
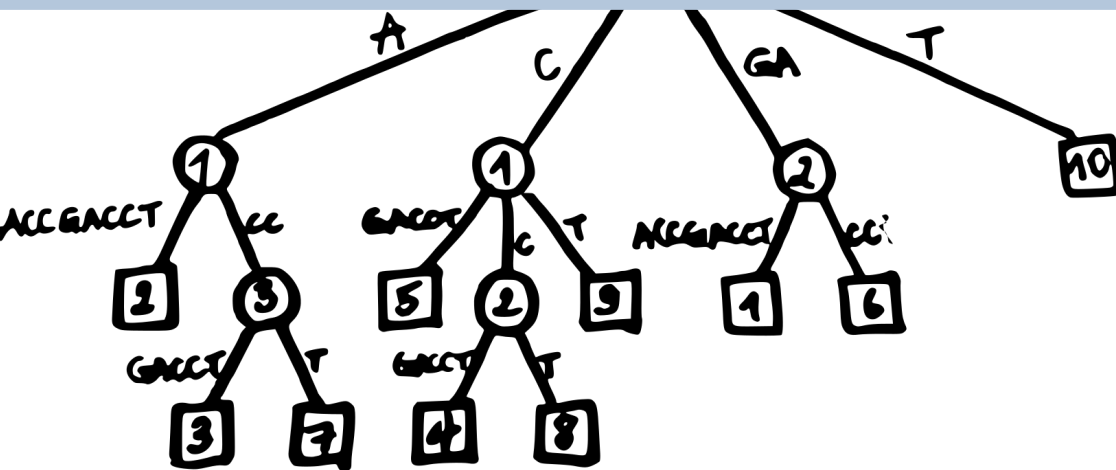
DAG workshop, 2020
Cyril Matthey-Doret

cyril.mattheydoret@gmail.com

cmdoret

# Exercises from session 1

# Exact sequence matching

Given two sequences (DNA, RNA, proteins, …)
At what position of target does query occur ?

query        CATAA

target  TAGACATAAA
          0123456789

          ↓

          4

When is this useful ?
- RNA quantification
- Read overlapping in genome assembly
- Primer design

# Exact sequence matching

Given two sequences (DNA, RNA, proteins, …)
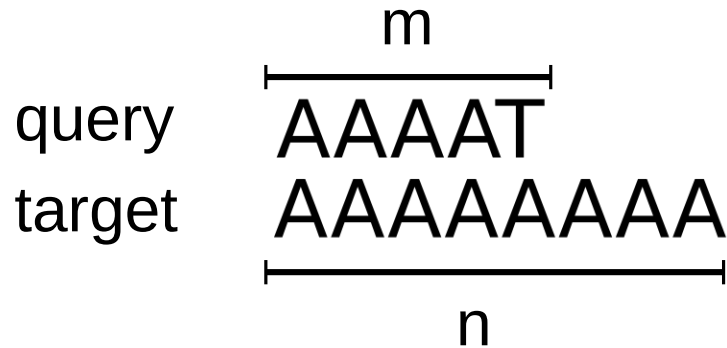At what position of target does query occur ?

query CATAA
target TAGACATAAA
0123456789

4

Can we come up with a naive algorithm ?

```
Pythonic way:
target = "TAGACATAAAGA"
query = "CATAA"
target.find(query)
4
```

# Naive solution

- How many possible alignments ?

- Number of comparisons ? (worst case)

m

query AAAAT
target AAAAAAAA

n

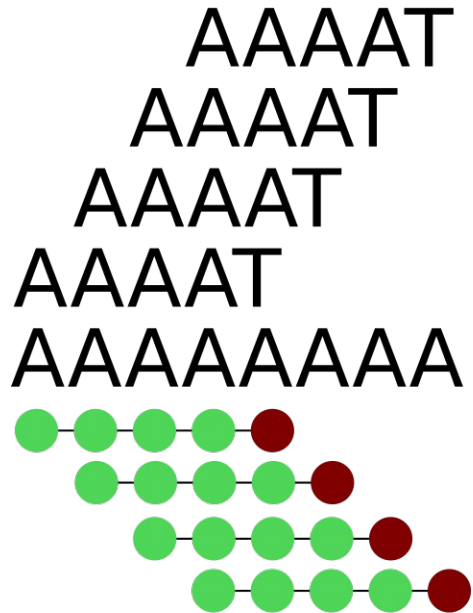# Naive solution

- How many possible alignments ? → n - m + 1

- Number of comparisons ? (worst case)

```
            AAAAT
             AAAAT
              AAAAT
query    AAAAT
target   AAAAAAAA
```

# Naive solution

- How many possible alignments ? → n - m + 1

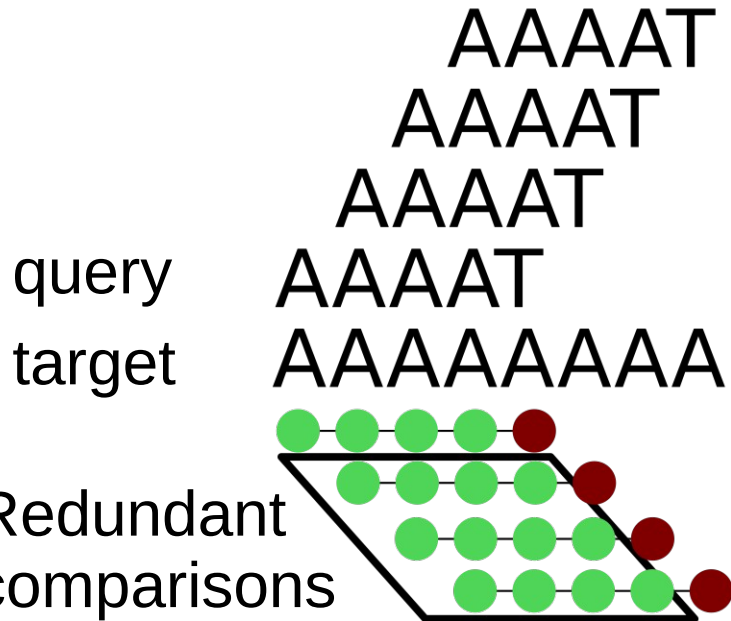- Number of comparisons ? (worst case) → m(n - m + 1)

AAAAT
AAAAT
AAAAT
query AAAAT
target AAAAAAAA

Big O notation ?

# Naive solution

- How many possible alignments ? → n - m + 1

- Number of comparisons ? (worst case) → m(n - m + 1)

AAAAT
AAAAT
AAAAT

query AAAAT

target AAAAAAAA

Redundant
comparisons

Big O notation ?
O(mn)

How to get faster ?

# The bad character rule

- Skip multiple positions at once

- One of the optimisations used in the Boyer-Moore algorithm

# The bad character rule

- Skip multiple positions at once

- One of the optimisations used in the Boyer-Moore algorithm

# The bad character rule

- Skip multiple positions at once

- One of the optimisations used in the Boyer-Moore algorithm

# The bad character rule

- Skip multiple positions at once

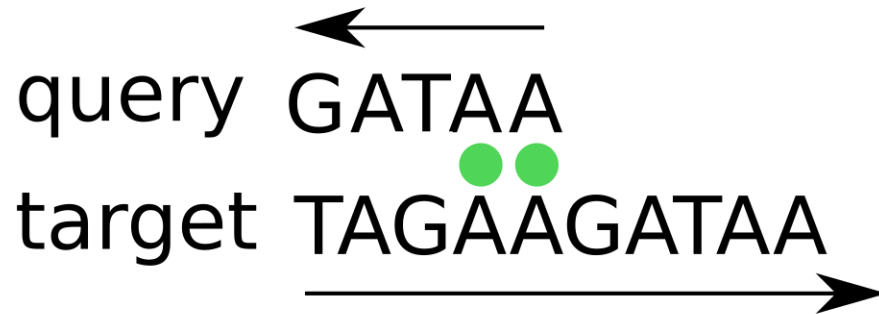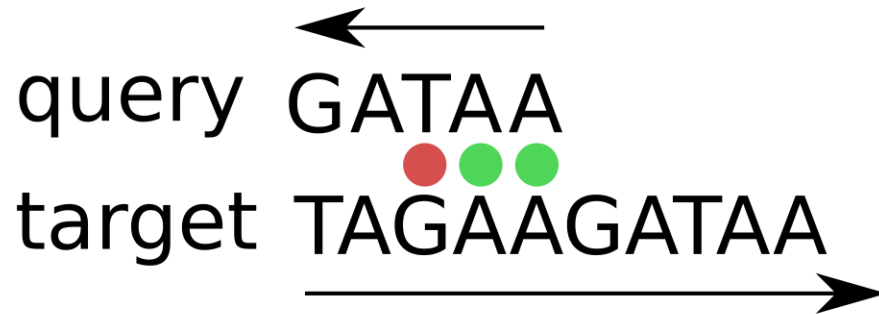- One of the optimisations used in the Boyer-Moore algorithm

# The bad character rule

- Skip multiple positions at once

- One of the optimisations used in the Boyer-Moore algorithm



Upon mismatch, skip until next matching character

# The bad character rule

- Skip multiple positions at once

- One of the optimisations used in the Boyer-Moore algorithm
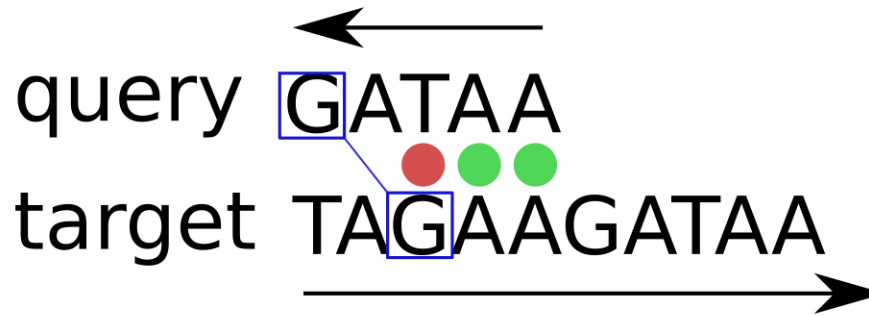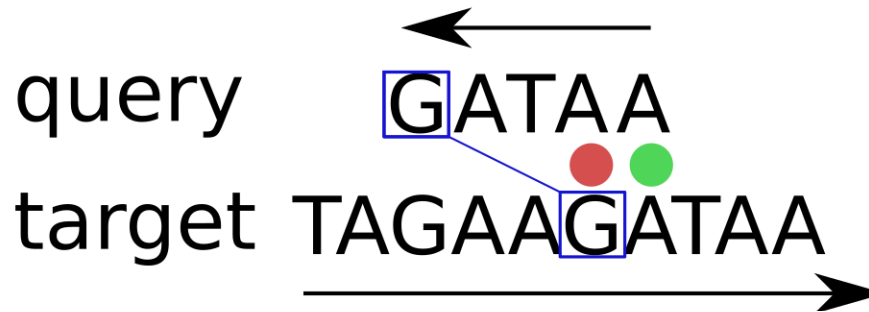
query     GATAA

target   TAGAAGATAA

# Preprocessing jumps

- Pre-compute jump sizes using a hash table !

- Scanning the query for the next match would take O(m)

- Looking up ["A-G"] in the hash table takes O(1)



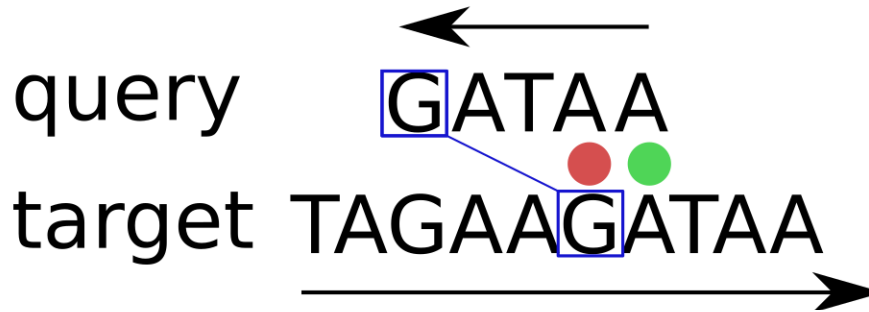|  | Query | | | | |
|---|---|---|---|---|---|
|  | G | A | T | A | A |
| A | 0 | - | 0 | - | - |
| C | 0 | 1 | 2 | 3 | 4 |
| T | 0 | 1 | - | 0 | 1 |
| G | - | 0 | 1 | 3 | 3 |

Alphabet

# Preprocessing jumps

- Pre-compute jump sizes using a hash table !

- Scanning the query for the next match would take O(m)

- Looking up ["A-G"] in the hash table takes O(1)



**This and other optimisations, lead to the Boyer-Moore algorithm, which is O(m+n)**

# Indexing

- Pre-processing the input sequence allows faster queries

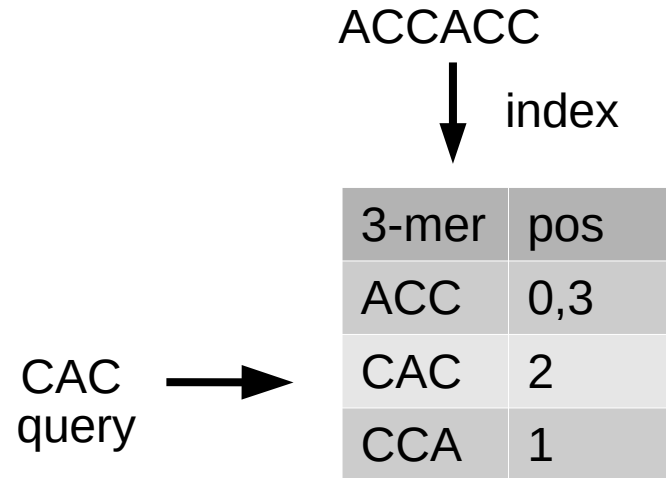- Different data structures can serve as an index

  - Hash table

  - suffix array

  - suffix tree

  - ...

When is this useful ?
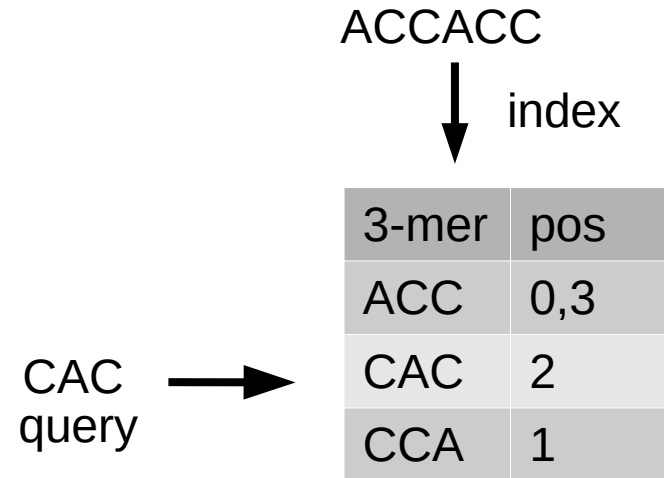- Read alignment
- Random access of large files

# Indexing

- Pre-processing the input sequence allows faster queries

- Different data structures can serve as an index

  - Hash table

  - suffix array

  - suffix tree

  - ...

ACCACC

↓ index

| 3-mer | pos |
|-------|-----|
| ACC   | 0,3 |
| CAC   | 2   |
| CCA   | 1   |

CAC
query →

# Indexing

- With a fix query size, we can use a k-mer based dictionary

- But what if the query size varies ?

ACCACC

↓ index

| 3-mer | pos |
|-------|-----|
| ACC   | 0,3 |
| CAC   | 2   |
| CCA   | 1   |

CAC
query →

# Indexing

- With a fix query size, we can use a k-mer based dictionary

- But what if the query size varies ?

  → Index all suffixes instead of k-mers

target  ATAGGGCA

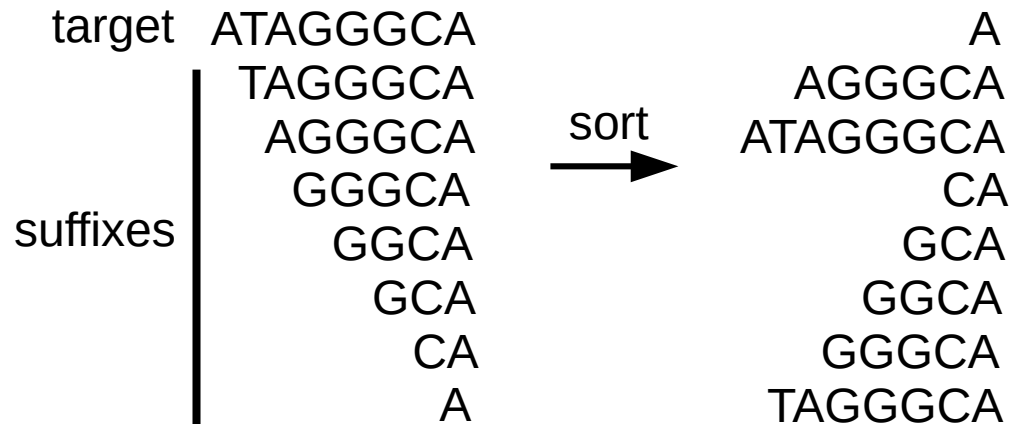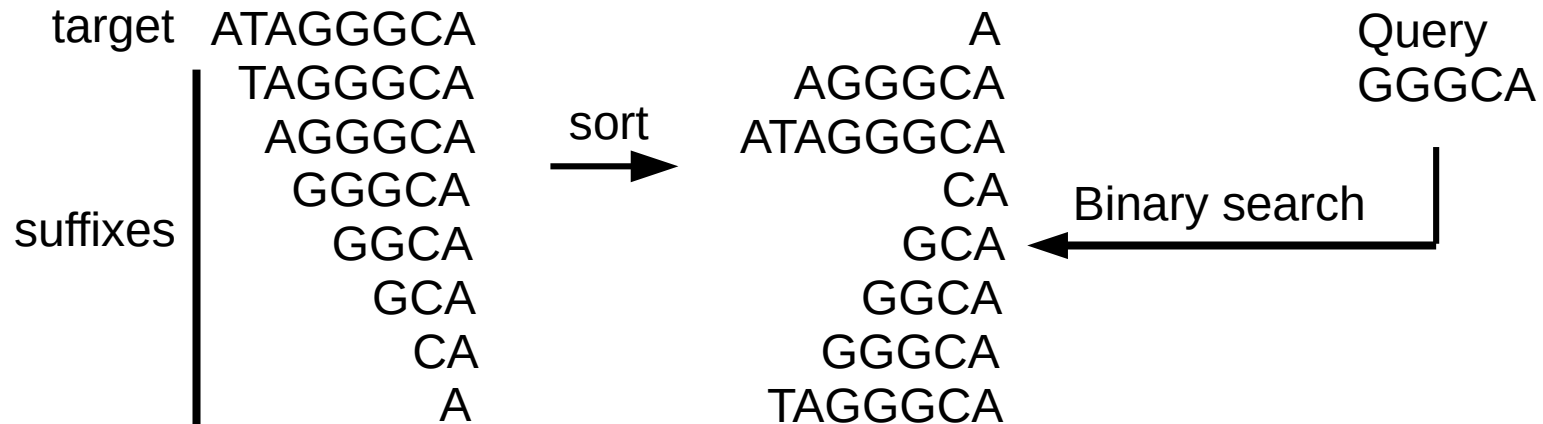TAGGGCA

AGGGCA

GGGCA

suffixes  GGCA

GCA

CA

A

# Indexing

- With a fix query size, we can use a k-mer based dictionary

- But what if the query size varies ?

    → Index all suffixes instead of k-mers

|        | target   ATAGGGCA |       |              A |
|--------|-------------------|-------|----------------|
|        | TAGGGCA           |       |         AGGGCA |
|        | AGGGCA            | sort  |       ATAGGGCA |
|        | GGGCA             | →     |             CA |
| suffixes | GGCA            |       |            GCA |
|        | GCA               |       |           GGCA |
|        | CA                |       |          GGGCA |
|        | A                 |       |        TAGGGCA |

# Indexing

- With a fix query size, we can use a k-mer based dictionary

- But what if the query size varies ?
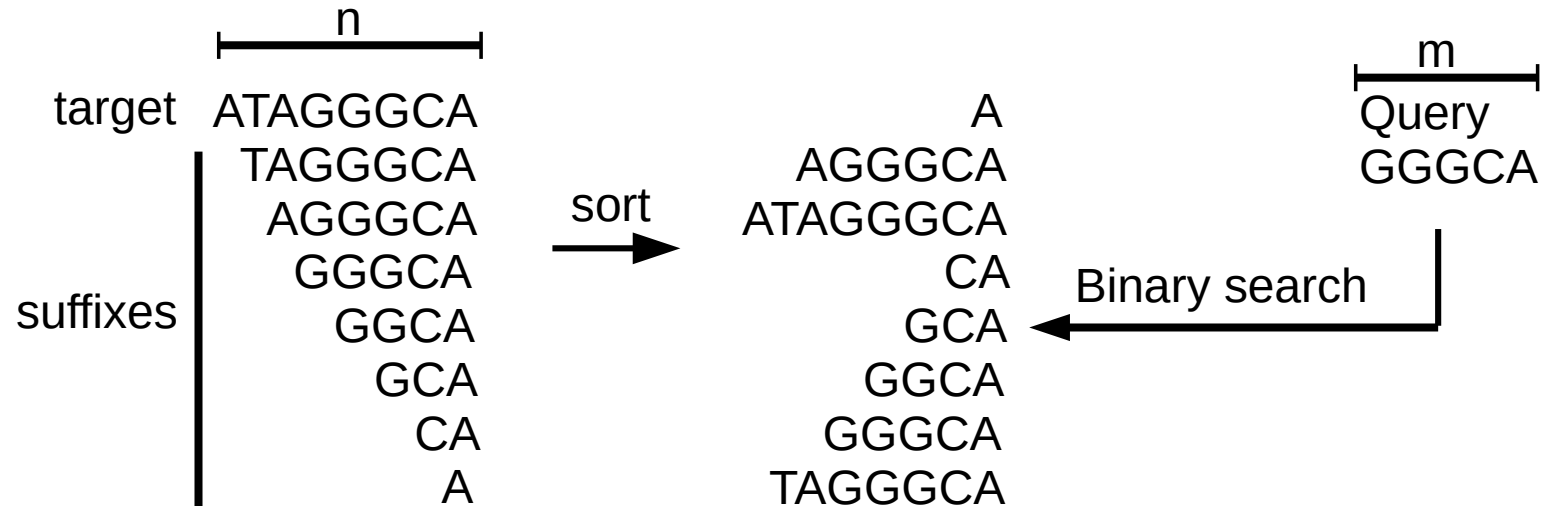
  → Index all suffixes instead of k-mers

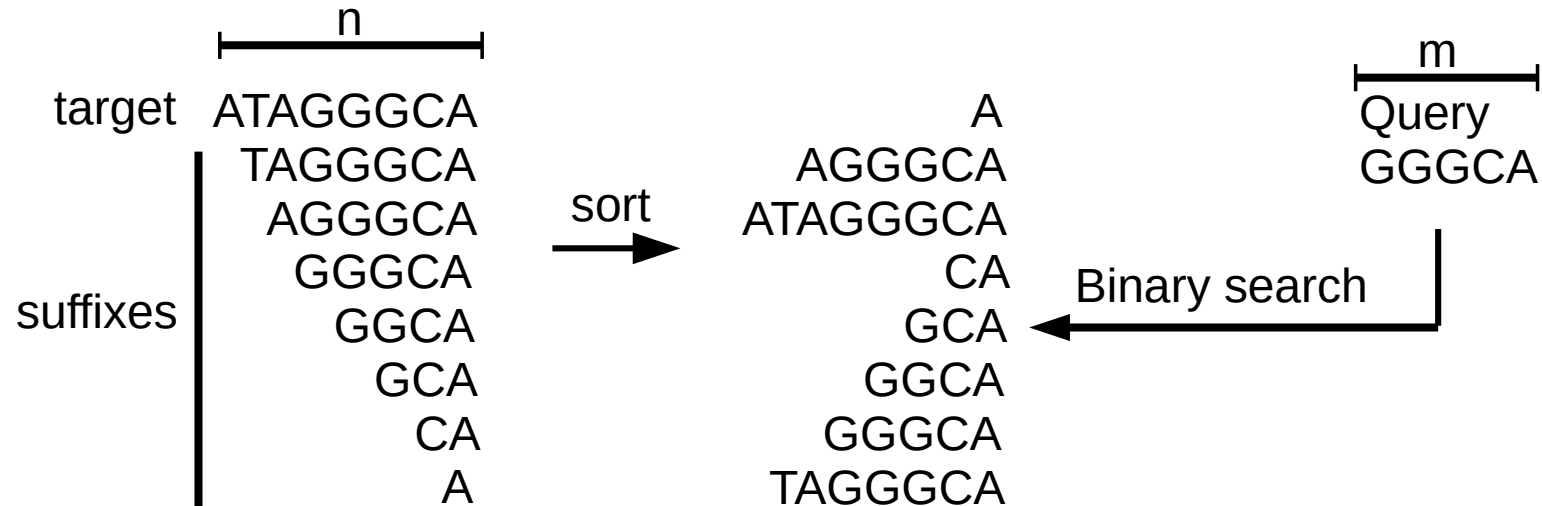| target | ATAGGGCA | | A | Query |
| | TAGGGCA | | AGGGCA | GGGCA |
| | AGGGCA | sort | ATAGGGCA | |
| | GGGCA | → | CA | Binary search |
| suffixes | GGCA | | GCA | ← |
| | GCA | | GGCA | |
| | CA | | GGGCA | |
| | A | | TAGGGCA | |

# Space analysis

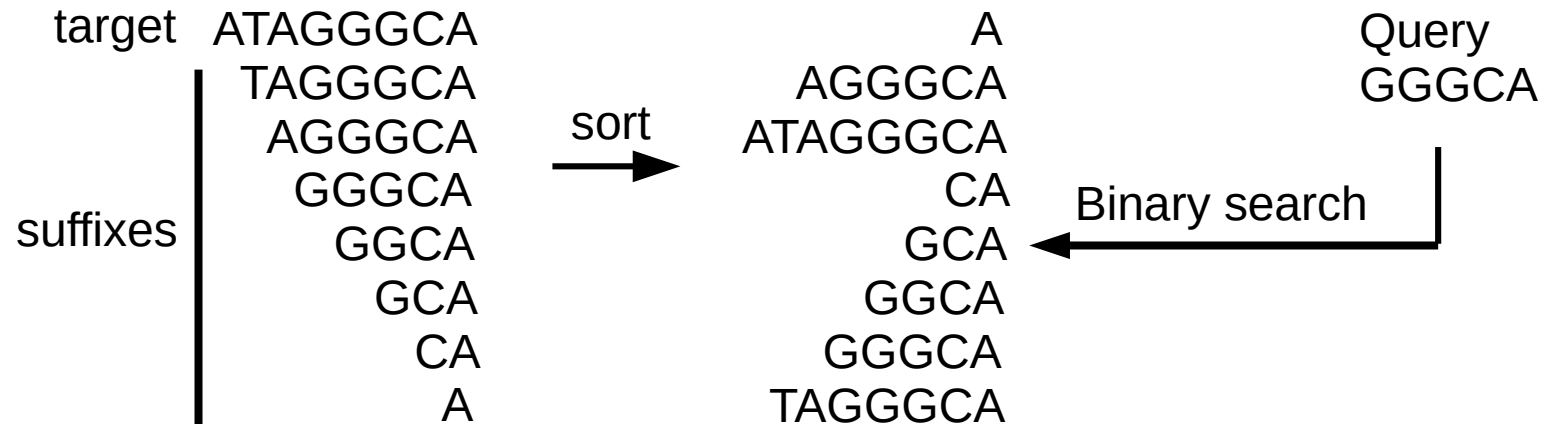- This is fast ! What is the query time complexity ?

# Space analysis

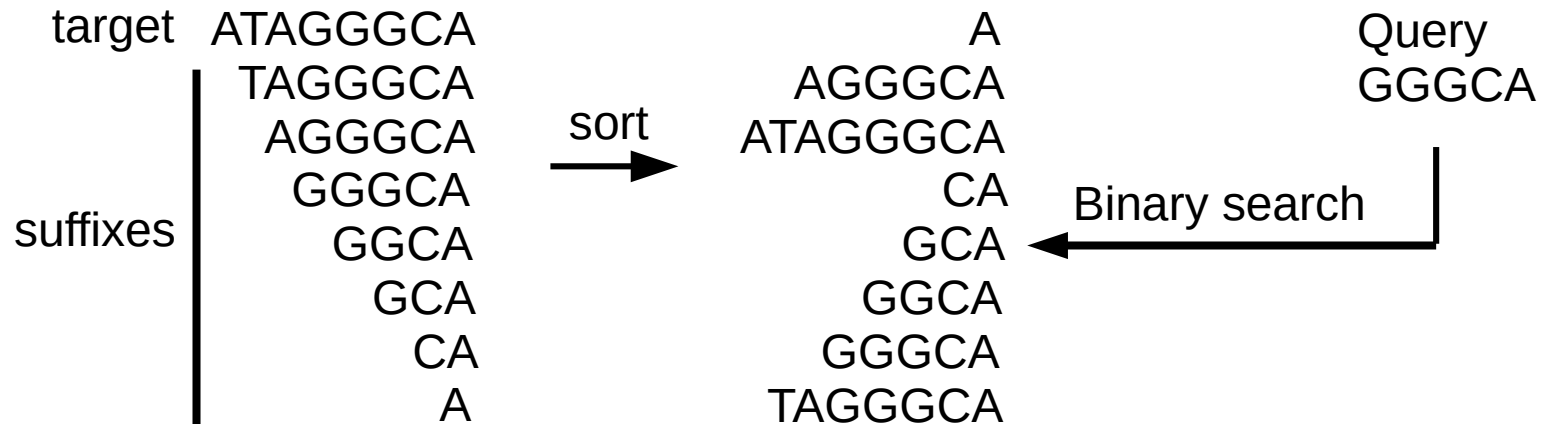- This is fast ! What is the query time complexity ? $O(m \log n)$

# Space analysis

- This is fast ! What is the query time complexity ? $O(m \log n)$

- But… It's huge ! How to estimate memory requirement ?

  → Big O notation

| target | ATAGGGCA | | | A | Query |
|---|---|---|---|---|---|
| | TAGGGCA | | | AGGGCA | GGGCA |
| | AGGGCA | sort | | ATAGGGCA | |
| | GGGCA | | | CA | |
| suffixes | GGCA | | | GCA | Binary search |
| | GCA | | | GGCA | |
| | CA | | | GGGCA | |
| | A | | | TAGGGCA | |

# Space analysis

- This is fast ! What is the query time complexity ? O(m log n)

- But… It's huge ! How to estimate memory requirement ?

$\rightarrow$ Big O notation: $n + (n - 1) + (n - 2) \ldots + 1 = \dfrac{n(n+1)}{2} = O(n^2)$

target  ATAGGGCA
TAGGGCA
AGGGCA
sort
GGGCA
suffixes  GGCA
GCA
CA
A

A
AGGGCA
ATAGGGCA
CA
GCA
GGCA
GGGCA
TAGGGCA

Query
GGGCA

Binary search

# Space analysis

Any ideas how to remove the redundant parts ?

```
        A
    AGGGCA
  ATAGGGCA
       CA
      GCA
     GGCA
    GGGCA
   TAGGGCA
```

# Space analysis

Any ideas how to remove the redundant parts ?
→ Only store the offsets

```
7          A
2      AGGGCA
0   ATAGGGCA
6         CA
5        GCA
4       GGCA
3      GGGCA
1   TAGGGCA
```
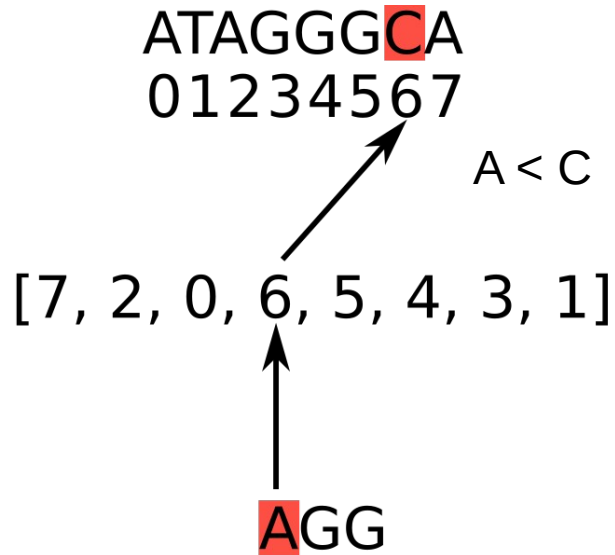
# Suffix array

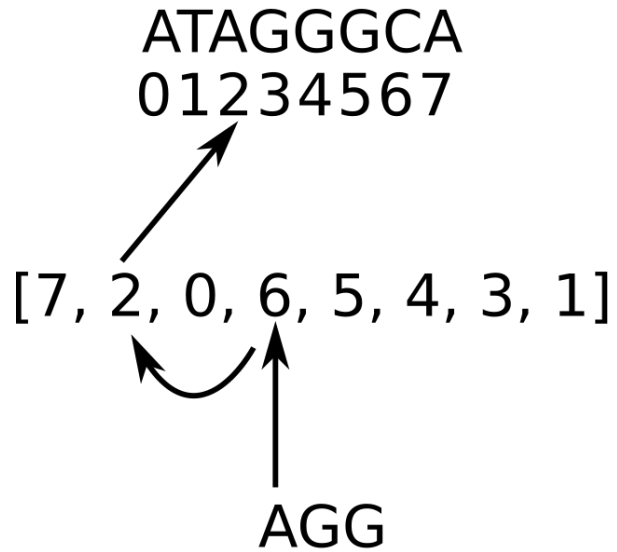- We can use binary search to look up the suffix (offset) array

ATAGGGCA
01234567

[7, 2, 0, 6, 5, 4, 3, 1]

AGG

# Suffix array

- We can use binary search to  look up the suffix (offset) array

ATAGGGCA
01234567

A < C

[7, 2, 0, 6, 5, 4, 3, 1]

AGG

# Suffix array

- We can use binary search to look up the suffix (offset) array

ATAGGGCA
01234567

[7, 2, 0, 6, 5, 4, 3, 1]

AGG

# Suffix array

- We can use binary search to look up the suffix (offset) array

ATAGGGCA
01234567

[7, 2, 0, 6, 5, 4, 3, 1]

AGG

# Suffix array

- We can use binary search to look up the suffix (offset) array

ATAGGGCA
01234567

[7, 2, 0, 6, 5, 4, 3, 1]

AGG

What is the search time for a pattern of length m ?

# Suffix array

- We can use binary search to  look up the suffix (offset) array

ATAGGGCA
01234567

[7, 2, 0, 6, 5, 4, 3, 1]

AGG

What is the search time for a pattern of

length m ? O(m log n)

# Suffix trees

- We can search in O(m) using a suffix tree
- Construction time is non-trivial (O(n))

T = "TACTCA$"
P = "CA$"

# Suffix trees

- We can search in O(m) using a suffix tree
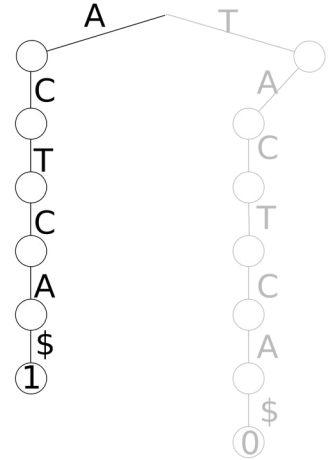- Construction time is non-trivial (O(n))

T = "TACTCA$"

# Suffix trees

- We can search in O(m) using a suffix tree
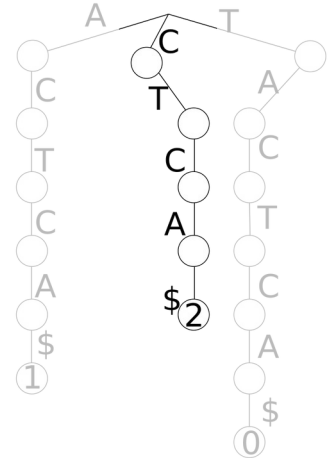- Construction time is non-trivial (O(n))

T = "TACTCA$"

# Suffix trees

- We can search in O(m) using a suffix tree
- Construction time is non-trivial (O(n))

T = "TACTCA$"

# Suffix trees

- We can search in O(m) using a suffix tree
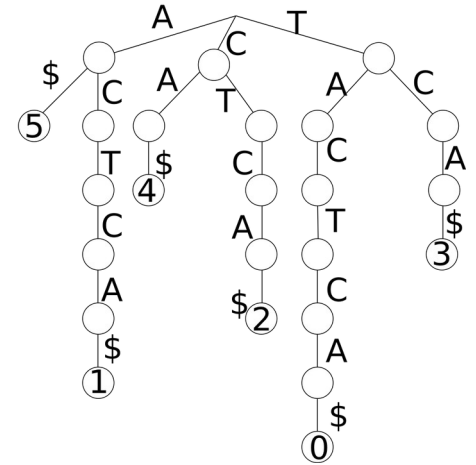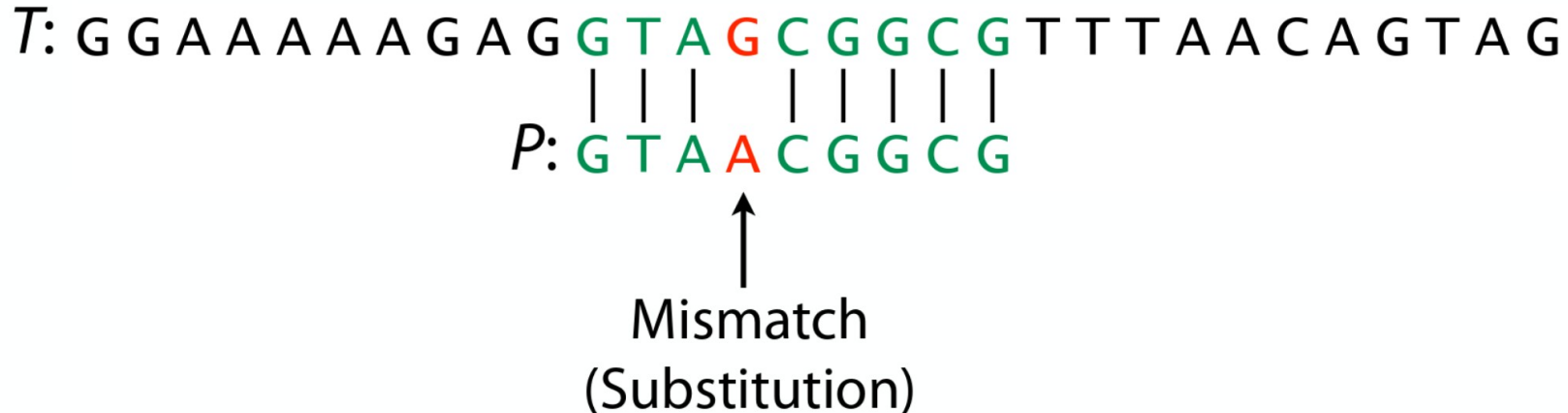- Construction time is non-trivial (O(n))

T = "TACTCA$"

# Inexact alignments

- We often want to allow *some* errors in alignments !
  - Sequencing errors in reads
  - Measure strain / species divergence
  - Find homologous genes
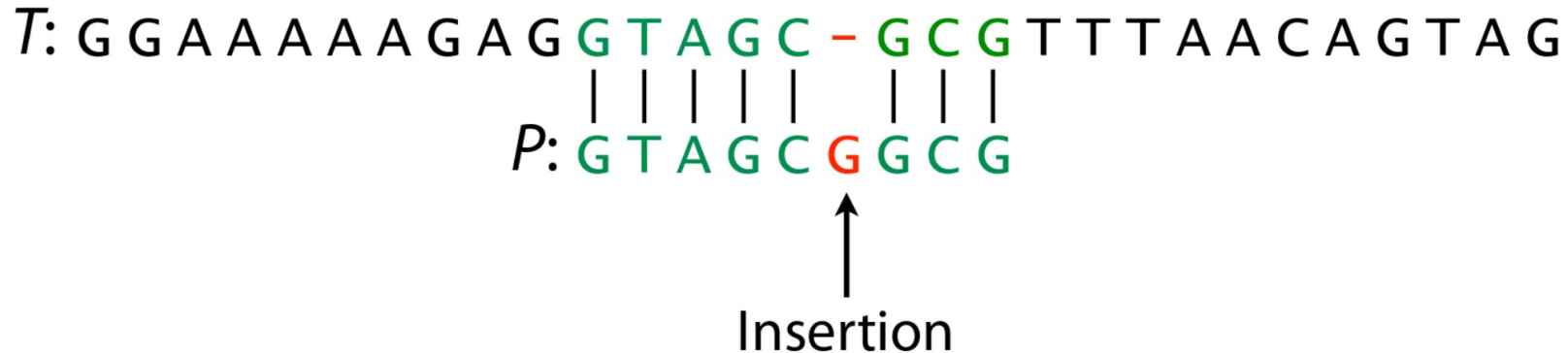


Mismatch
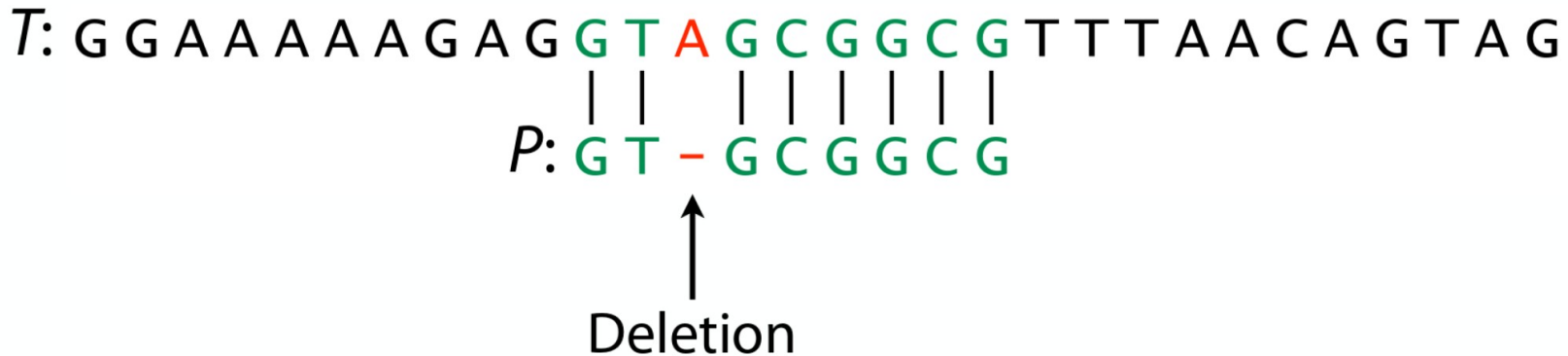(Substitution)

# Inexact alignments

- We often want to allow *some* errors in alignments !
  - Sequencing errors in reads
  - Measure strain / species divergence
  - Find homologous genes

# Inexact alignments

- We often want to allow *some* errors in alignments !
  - Sequencing errors in reads
  - Measure strain / species divergence
  - Find homologous genes


Deletion

# Quantifying errors

- Hamming distance: # of substitutions to turn P into T
- Edit distance: Same, but allows insertions / deletions

T = TTGCC
P = CTCGC

⬇

TTGCC
TTCGC

Substitute C → T

⬇

TTGCC
TTGGC

Substitute C → G

What if P ant T have different lengths ?

Hamming distance = 2

# Inexact alignments

- Aligned sequence will almost always have different lengths
- This can be solved with 2 approaches

$$T = \overline{TTGCCC}^{\,n}$$
$$P = \underline{TGCG}_{\,m}$$

**Global alignment**

Naive: compute all possible alignments

```
TTGCCC      TTGCCC      TTGCCC
TGCG__      TGC_G_      TGC__G

TTGCCC      TTGCCC      ...
TC_C_G      _TGC_G
```

**Local alignment**

Search for substring of length m in T with smallest edit distance

```
TTGCCC
 TGCG
```

# Inexact alignments

- Aligned sequence will almost always have different lengths
- This can be solved with 2 approaches

$$\overline{\phantom{TTGCCC}} \; n$$

$$T = TTGCCC$$
$$P = \underline{TGCG}$$

$$m$$

Global alignment                                                                 Local alignment

Edit distance: Too many possibilities !

# Inexact alignments

- Can be formulated as a recursive problem
- Large exploration space

d(AACCTG, CCCG) =

$$
min
\begin{cases}
d(AACCT, CCCG) + 1 = 3 + 1 = 4 \\
\\
d(AACCT, CCC) \;\; + 0 = 3 + 0 = \boxed{3} \\
\\
d(AACCTG, CCC) + 1 = 4 + 1 = 5
\end{cases}
$$

AACCT G
CCCG _

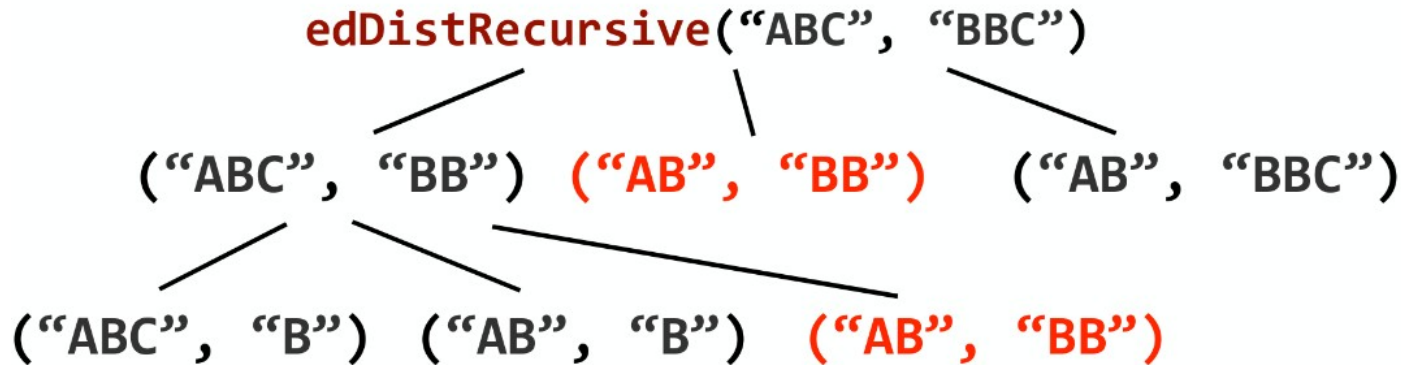*assume this part solved*

AACCT G
CCC G

AACCTG _
CCC G

# Inexact alignments

- Can be formulated as a recursive problem
- How to reduce the exploration space ?

```python
def edDistRecursive(a, b):
    if len(a) == 0:
        return len(b)
    if len(b) == 0:
        return len(a)
    delt = 1 if a[-1] != b[-1] else 0
    return min(edDistRecursive(a[:-1], b[:-1]) + delt,
               edDistRecursive(a[:-1], b) + 1,
               edDistRecursive(a, b[:-1]) + 1)
```

# Inexact alignment in feasible time

- Can be formulated as a recursive problem
- How to reduce the exploration space ?

# Dynamic programming

- We perform "traceback to recover the alignment from scores

| | (empty string) | **A** | A**A** | AA**C** | AAC**C** | AACC**T** | AACCT**G** |
|---|---|---|---|---|---|---|---|
| (empty string) | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| **C** | 1 | 1 | 2 | 2 | 3 | 4 | 5 |
| C**C** | 2 | 2 | 2 | 2 | 2 | 3 | 4 |
| CC**C** | 3 | 3 | 3 | 2 | 2 | 3 | 4 |
| CCC**G** | 4 | 4 | 4 | 3 | 3 | 3 | 3 |

Recursion equation

$$D_{i,j} = \min \begin{cases} D_{i-1,j-1} & + & -1 & a_i = b_j \\ D_{i-1,j-1} & + & 1 & a_i \neq b_j \\ D_{i-1,j} & + & 1 & b_j = - \\ D_{i,j-1} & + & 1 & a_i = - \end{cases}$$

AACCTG
__CCCG

# Dynamic programming

- We perform "traceback to recover the alignment from scores



Recursion equation

$$D_{i,j} = \min \begin{cases} D_{i-1,j-1} & + & -1 & a_i = b_j \\ D_{i-1,j-1} & + & 1 & a_i \neq b_j \\ D_{i-1,j} & + & 1 & b_j = - \\ D_{i,j-1} & + & 1 & a_i = - \end{cases}$$
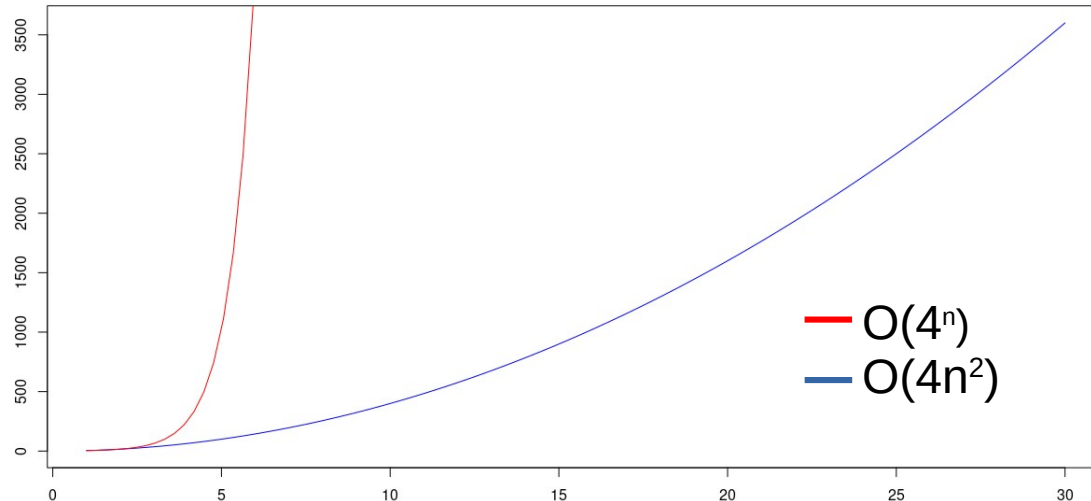
AACCTG
__CCCG

**Live visualisation here !**

# Dynamic programming

- For any pair of sequence prefixes, score is calculated once
- Time complexity ?
  - how many operations /cell ?
  - How many cells ?

# Dynamic programming

- For any pair of sequence prefixes, score is calculated once
- Time complexity ? O(mn)
  - how many operations /cell ? **3 addition and 1 minimum**
  - How many cells ? **(m + 1) (n + 1)**

# Refining dynamic programming

- Add different penalties (e.g. match=1, gap=-2, mismatch=-1)
- Use a substitution matrix (measures amino-acid similarity)



A→ R: -3
A → W: -7

# Local alignment

- Smith-waterman uses Dynamic programming for local alignments

  - All negative scores set to 0

  - Stop extending when encountering zeros

  - Start backtracking from highest score

**Initialize the scoring matrix**

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

Substitution matrix: $S(a_i, b_j) = \begin{cases} +3, & a_i = b_j \\ -3, & a_i \neq b_j \end{cases}$

Gap penalty: $W_k = kW_1$
$W_1 = 2$

Gif: https://w.wiki/ZAL

# Dynamic programming in real life

- Used in most alignment tools

- Often combined with a first step of exact matching for speed

- Example: BLAST

    1. Extract k-mers from the query

    2. Find their position in the database (exact match)

    3. Identify regions with several exact matches as "High scoring pairs"

    4. Extend HSPs in both directions using dynamic programming

# Additional resources

- **Book on string-based algorithms**: Algorithms on Strings, Trees, and Sequences Computer Science and Computational Biology, Dan Gusfield