# Computational Biology
## DNA Sequencing

Department of Mathematics
Stockholm University
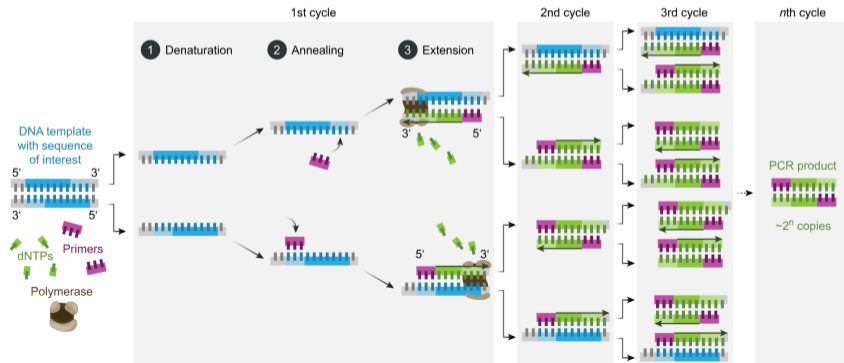
- Copying DNA:
  - Polymerase chain reaction (PCR)
- Sequencing DNA:
  - Sanger Sequencing [AKA 1st generation sequencing]
  - Next/2nd-generation sequencing (NGS) [AKA Massive parallel sequencing]
  - 3rd-generation [AKA long-read sequencing]

## Polymerase chain reaction (PCR)

- used to copy DNA
- Invented by Kary Mullis (Nobel prize 1993)
- Input: a DNA "template" $t$ to copy, primers, polymerase, bases $A, C, G, T$,

  Process: $n$ "cycles" (see right)

  Output: roughly $2^n$ copies of $t$

Per cycle there are 3 phases:

1. Denaturate: 94-98 °C for 20–30 s
2. Anneal: 50-65 °C for 20–40 s
3. Extension: 75-80 °C



fig taken from https://en.wikipedia.org/wiki/Polymerase_chain_reaction

## Sanger Sequencing
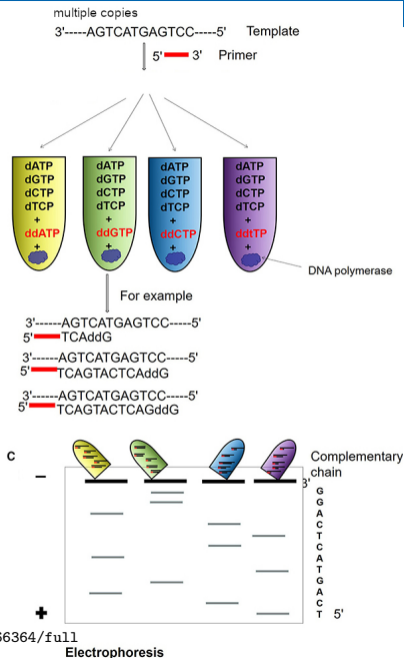
▶ used to read "small ( 500bp)" DNA sequences

▶ Invented by Fredrick Sanger and coworkers, 1977 (Nobel prize 1980)

▶ Input: copies of DNA split into 4 test tubes that contains primers, polmerase, bases, "modified bases $A, C, T, G$"

Each tube contains all bases and ONE "modified base" $I \in \{A, C, G, T\}$

Process (Basic Idea): "modified base" $I$ ensures that when added during reading process of one DNA-copy, the reading process stops.

Having multiple copies and the four tubes, this ensures, that (with high probability) the tupe $I$ contains all single strands that end with $I$.

gel electrophoresis: reads are negative charged and small reads get "closer" to positive pol (proportional to their length)

Output: the read of the input DNA

fig taken from https://www.frontiersin.org/articles/10.3389/fmicb.2021.766364/full
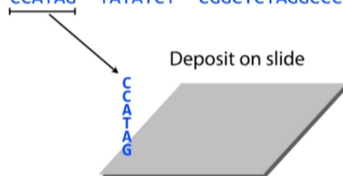


Electrophoresis

## Next-generation sequencing (NGS)

- ▶ used to read **multiple** "small ( 500bp)" DNA sequences

- ▶ Several methods exits, one is the "Illumina sequencing process":

  Input: copies of **multiple** DNA (fragments) placed on a slide, bases, terminators, polymerase, ..



figs taken from Ben Langmead "ADS1: Sequencing by Synthesis" (youtube)

## Next-generation sequencing (NGS)

- ▶ used to read **multiple** "small ( 500bp)" DNA sequences

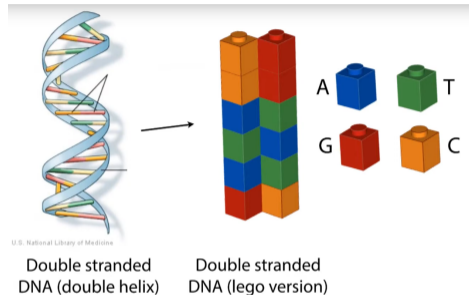- ▶ Several methods exits, one is the "Illumina sequencing process":

  Input: copies of **multiple** DNA (fragments) placed on a slide, bases, terminators, polymerase, ..



Double stranded DNA (double helix)

Double stranded DNA (lego version)

figs taken from Ben Langmead "ADS1: Sequencing by Synthesis" (youtube)
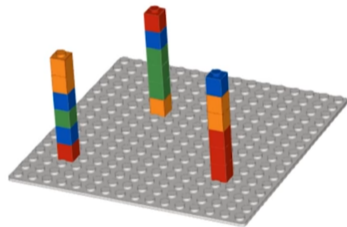
## Next-generation sequencing (NGS)

▶ used to read **multiple** "small ( 500bp)" DNA sequences

▶ Several methods exits, one is the "Illumina sequencing process":
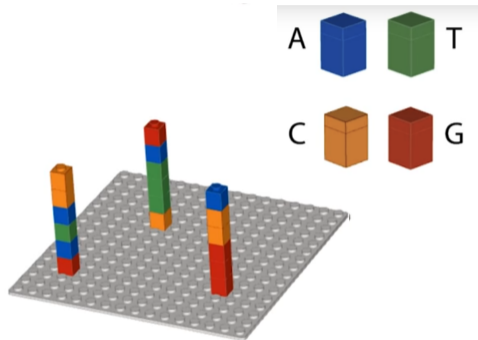
Input: copies of **multiple** DNA (fragments) placed on a slide, bases, terminators, polymerase, ..

# Next-generation sequencing (NGS)

- ▶ used to read **multiple** "small ( 500bp)" DNA sequences

- ▶ Several methods exits, one is the "Illumina sequencing process":

  Input: copies of **multiple** DNA (fragments) placed on a slide, bases, terminators, polymerase, ..



figs taken from Ben Langmead "ADS1: Sequencing by Synthesis" (youtube)

## Next-generation sequencing (NGS)

▶ used to read **multiple** "small ( 500bp)" DNA sequences

▶ Several methods exits, one is the "Illumina sequencing process":

Input: copies of **multiple** DNA (fragments) placed on a slide, bases, terminators, polymerase, ..

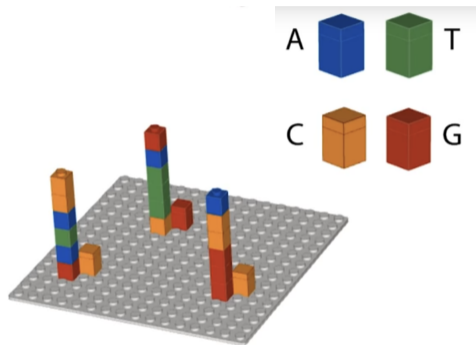Process (Basic Idea): when bases with terminators bind, no further base can be added.

# Next-generation sequencing (NGS)

- used to read **multiple** "small ( 500bp)" DNA sequences

- Several methods exits, one is the "Illumina sequencing process":

  Input: copies of **multiple** DNA (fragments) placed on a slide, bases, terminators, polymerase, ..

  Process (Basic Idea): when bases with terminators bind, no further base can be added.

  terminators are engineered to glow a particular color ($A, C, G, T$)
  $\rightarrow$ take photo



figs taken from Ben Langmead "ADS1: Sequencing by Synthesis" (youtube)
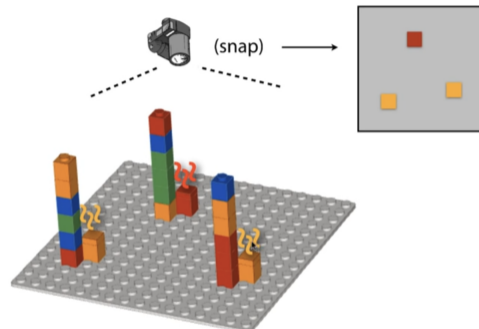
## Next-generation sequencing (NGS)

▶ used to read **multiple** "small ( 500bp)" DNA sequences

▶ Several methods exits, one is the "Illumina sequencing process":

Input: copies of **multiple** DNA (fragments) placed on a slide, bases, terminators, polymerase, ..

Process (Basic Idea): when bases with terminators bind, no further base can be added.

terminators are engineered to glow a particular color ($A, C, G, T$)
$\rightarrow$ take photo

after taking photo, terminators are removed and process is repeated.



Remove terminators

---

## Next-generation sequencing (NGS)

- ▶ used to read **multiple** "small ( 500bp)" DNA sequences

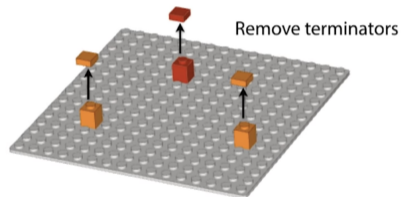- ▶ Several methods exits, one is the "Illumina sequencing process":

  Input: copies of **multiple** DNA (fragments) placed on a slide, bases, terminators, polymerase, ..

  Process (Basic Idea): when bases with terminators bind, no further base can be added.

  terminators are engineered to glow a particular color ($A$, $C$, $G$, $T$)
  $\rightarrow$ take photo

  after taking photo, terminators are removed and process is repeated.



---

figs taken from Ben Langmead "ADS1: Sequencing by Synthesis" (youtube)

## Next-generation sequencing (NGS)

▶ used to read **multiple** "small ( 500bp)" DNA sequences

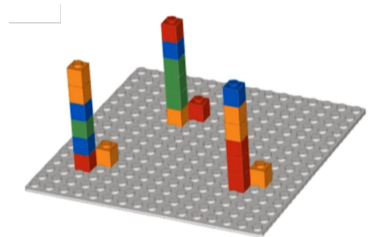▶ Several methods exits, one is the "Illumina sequencing process":

Input: copies of **multiple** DNA (fragments) placed on a slide, bases, terminators, polymerase, ..

Process (Basic Idea): when bases with terminators bind, no further base can be added.

terminators are engineered to glow a particular color ($A$, $C$, $G$, $T$)
$\rightarrow$ take photo

after taking photo, terminators are removed and process is repeated.



figs taken from Ben Langmead "ADS1: Sequencing by Synthesis" (youtube)

## Next-generation sequencing (NGS)

▶ used to read **multiple** "small ( 500bp)" DNA sequences

▶ Several methods exits, one is the "Illumina sequencing process":

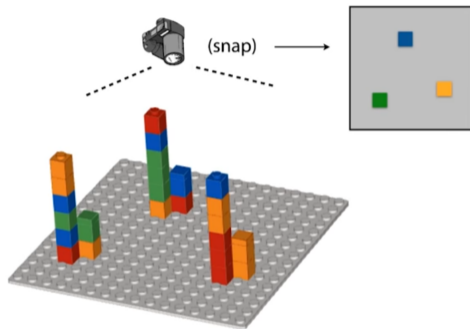Input: copies of **multiple** DNA (fragments) placed on a slide, bases, terminators, polymerase, ..

Process (Basic Idea): when bases with terminators bind, no further base can be added.

terminators are engineered to glow a particular color ($A$, $C$, $G$, $T$)
$\rightarrow$ take photo

after taking photo, terminators are removed and process is repeated.



figs taken from Ben Langmead "ADS1: Sequencing by Synthesis" (youtube)

## Next-generation sequencing (NGS)

- ▶ used to read **multiple** "small ( 500bp)" DNA sequences

- ▶ Several methods exits, one is the "Illumina sequencing process":

  Input: copies of **multiple** DNA (fragments) placed on a slide, bases, terminators, polymerase, ..

  Process (Basic Idea): when bases with terminators bind, no further base can be added.

  terminators are engineered to glow a particular color ($A$, $C$, $G$, $T$)
  → take photo

  after taking photo, terminators are removed and process is repeated.

  Output: the read of the **multiple** input DNAs (photos of each cycle)



figs taken from Ben Langmead "ADS1: Sequencing by Synthesis" (youtube)

## Next-generation sequencing (NGS)

- ▶ used to read **multiple** "small ( 500bp)" DNA sequences

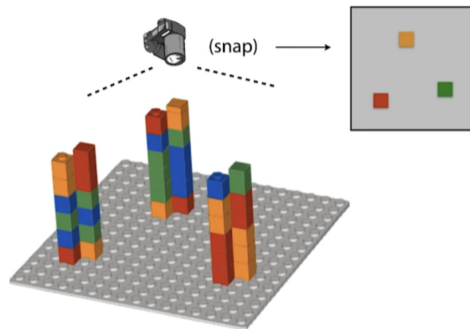- ▶ Several methods exits, one is the "Illumina sequencing process":

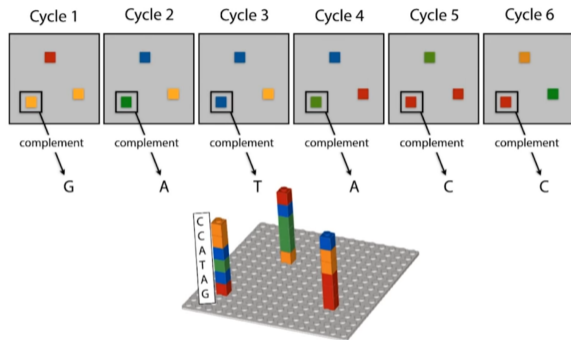  Input: copies of **multiple** DNA (fragments) placed on a slide, bases, terminators, polymerase, ..

  Process (Basic Idea): when bases with terminators bind, no further base can be added.

  terminators are engineered to glow a particular color ($A, C, G, T$)
  $\rightarrow$ take photo

  after taking photo, terminators are removed and process is repeated.

  Output: the read of the **multiple** input DNAs (photos of each cycle)



**Key feature:**
massively parallel, photograph captures all templates simultaneously (billions of DNA templates on a slide)

fig taken from https://theory.labster.com/ngs-experiment/

Cost per Human Genome

- Copying DNA:
  - Polymerase chain reaction (PCR)
- Sequencing DNA:
  - Sanger Sequencing [AKA 1st generation sequencing]
  - Next/2nd-generation sequencing (NGS) [AKA Massive parallel sequencing]
  - 3rd-generation [AKA long-read sequencing]
    (currently under active development[*], can read more than 10000 bp)

To recall, humanDNA $3.2 \times 10^9$bp, Carsonella ruddii DNA 159 662bp

Observation: Whole genomes cannot be read at once.

---

[*]Marx, V. Method of the year: long-read sequencing. Nat Methods 20, 6–11 (2023).

unknown DNA
??????????????????????????

**Observation:**

Sequencers cannot read whole genomes at once.

Idea_1: randomly break-up long DNA into multiple pieces
  (e.g. with ultrasound)
  and sequence them

However: if we just use a single DNA strand, well ..

break-up into random
small pieces

```
  ?????   ????   ?????   ???
  ?????   ????   ?????   ???
```

sequence
small pieces

```
  GTTAG           TGCAT   AAA
  AGGCT   CATT           GGC
```

assembly

```
TGCATGTTAGCATTAGGCTAAAGGC
GGCGTTAGTGCATGGCTCATTAAA
CATTTGCATGTTCAGGCTAAAGGC
```

**Observation:**

Sequencers cannot read whole genomes at once.

**Idea_1:** randomly break-up long DNA into multiple pieces

(e.g. with ultrasound)

and sequence them

However: if we just use a single DNA strand, well ..

**Idea_2:** Produce multiple copies of DNA first and then apply Idea_1

$\Longrightarrow$ results in overlapping reads

$\Longrightarrow$ assemply (here smart computational methods are needed!)

# Sequence Assembly

For a given set $\zeta = \{S_1, \ldots S_N\}$ of strings (=reads of fragments of DNA $D$),
a superstring is a string $S$ that contains all $S_i$ as substrings.

Trivially, we could concatenate all strings in $\zeta$ to get superstring $S$. However, having say $\sim 10^6$
copies of DNA $D$ fragmented and sequenced, we get then a string $S$ of length $|S| \sim |D| \times 10^6$
$\implies$ far away from $D$.

In the assembly problem, we want to find a superstring that "best represents" $D$.

There are several ways on how to define "best represents" !!

We start with considering following problem:

**Shortest Common Superstring Problem (SCS):**
For a given $\zeta = \{S_1, \ldots S_N\}$ find a superstring $S$ of shortest length.

SCS is NP-hard. So we focus ways to approximate solutions

$\implies$ overlap graphs and `Greedy_SCS` (board)

# Basics:

S: 
$$1\ 2\ 3\ 4\ 5\ 6\ 7\ 8$$
$$\text{A T T C G T A C}$$

$S[1..4]$ prefix

$S[7..8]$ suffix

$S(5) = L$

string $S = x_1 \cdots x_m$, $|S| = m$

$S[i..j] = x_i x_{i+1} \cdots x_j$

$S(i) = x_i$

$S[1..j] \hat{=}$ prefix of $S$ ending at $j$

$S[j..m] \hat{=}$ suffix of $S$ starting at $j$

$S'$ __substring of__ $S = x_1 \cdots x_m$, in symbols $S' \sqsubseteq S$

if $S' = x_i x_{i+1} \cdots x_j$, $1 \leq i < j \leq m$

Shotgun ~ assembly

Formally:      given DNA $D \to$ copy, i.e. we get

set $\{ D_1 \cdots D_e \}$ of copies of $D$

$\to$

Shotgun sequencing:

copy (

unknown DNA $D$

$D_1$
⋮
$D_e$

shotgun (

$F_1$
⋮
$F_N$   fragments

sequence (

$S_1$   know DNA fragments
⋮
$S_N$

assemble

reconstructed DNA $S$

we want to understand this step

## Assembly Problem:

For given set $S = \{S_1 \ldots S_N\}$ of strings (substrings of $D$)
find superstring $S$ s.t. $S_i \subseteq S$, $1 \leq i \leq N$
and $S$ "best represents" $D$.

There are several ways to define "best represents"
here we focus on:

## Shortest common superstring problem (SCS):

Given $S = \{S_1 \ldots S_N\}$. Find $S$ of min-length
s.t. $S_i \subseteq S$, $1 \leq i \leq N$.

EXMPL: $S = \{$ ATAT , TATA,
TATT , TAAT,
TTAT , AATA $\}$

// in this example,
$f_i$ of same length
but this is not necessary
for alg. //

$\longrightarrow$ shortest superstring $S$ is

T A A T A T T A T A

__Theorem (Gallant et al 1980):__   SCS is NP-hard

→ no hope for polynomial-time exact algorithm $(P \neq NP)$

Q: Can we approximate a solution?

__In what follows:__   $\mathcal{S} = \{S_1 .. S_N\}$   $\swarrow$ "substring-free"

$\supset t$   $S_i \neq S_j$ & $S_i \not\sqsubseteq S_j$ $\forall i \neq j$:

$\left[\text{IF } S_i = S_j \text{ or } S_i \sqsubseteq S_j , \text{ we can remove } S_i \atop \text{in preprocess step}\right]$

has no impact on SCS!

we will see later,
this can be done in
$O(N \cdot \|\mathcal{S}\|)$ time

$\|\mathcal{S}\| = \sum\limits_{S \in \mathcal{S}} |S|$   via Suffix-trees

[Lit: Algorithmic Aspects
of Bioinformatics, Böckenhauer & Bongartz,
2007, Springer]

**DEF:**

For all $S_i, S_j \in \mathcal{S} = \{S_1 \dots S_N\}$ there is
a longest substring $v_{ij}$ of $v_{ij}$ is suffix of $S_i$
& prefix of $S_j$

$$\text{that is} \quad S_i = u\, v_{ij}$$
$$S_j = v_{ij}\, w$$

$(v_{ij} = \varepsilon$ empty string
is possible! $)$

$v_{ij}$ called $\underline{\text{ovulap}}$ of $S_i$ & $S_j$

$$ov(i,j) = |v_{ij}| \qquad (!\ v_{ij} \neq v_{ji} \text{ possible!})$$
$$p(i,j) = |u|$$
$$= |S_i| - ov(i,j)$$
$$\text{merge}(i,j) = u\, v_{ij}\, w$$
$$\text{pref}(i,j) = u$$

**EXMPL** $\quad S_1 = \text{AT AT} \quad , \quad S_2 = \text{ATTT}$

$V_{12} = \text{AT}$

$ov(1,2) = 2$

$p(1,2) = 2$

$\text{merge}(1,2) = u\, v_{12}\, w$
$\qquad\qquad\quad = \text{ATATTT}$

$S_2 = \text{ATTT} \quad S_1 = \text{AT AT} \quad V_{21} = \varepsilon \text{ (empty)}$

$ov(2,1) = 0$

$p(2,1) = 4 \qquad [\, f_2 = \text{ATTT} = u\, v_{21} \,]$

$\text{merge}(2,1) = \text{ATTTATAT}$

**DEF:** graph $G_\mathcal{S} = (V, E)$ over $\mathcal{S} = \{S_1 \dots S_N\}$:

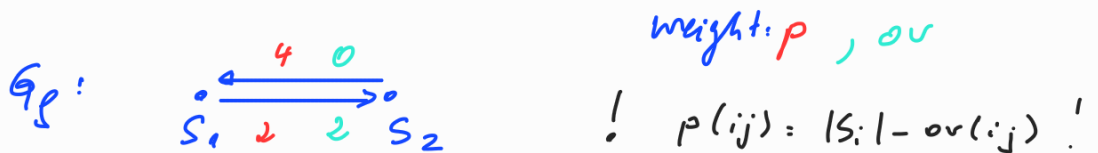- $V = \mathcal{S}$
- $E = (\mathcal{S} \times \mathcal{S}) \setminus \{(S,S) : S \in \mathcal{S}\}$

"no loops"

"prefix graph" $=:$ arc $S_i \longrightarrow S_j$ gets weight $p(i,j)$

"overlap graph": $\underline{\quad}$ " $\underline{\quad}$ $ov(i,j)$.

[can be constructed in $O(n \|\mathcal{S}\|)$ time]
($\nearrow$ suffixtrees later)

A <u>path</u> in $G_\mathcal{S}$ is a subgraph of the "form"

$$v_1 \longrightarrow v_2 \longrightarrow \dots \longrightarrow v_\ell$$

$$v_i \neq v_j \ \forall \ i \neq j,$$

**EXMPL:** (1) $\mathcal{S} = \{S_1 = ATAT, \ S_2 = ATTT\}$

$G_\mathcal{S}$:

weight: $p$, $ov$

$$\overset{4 \quad 0}{\underset{S_1 \ 2 \quad 2 \ S_2}{\rightleftarrows}}$$

! $p(ij) = |S_i| - ov(ij)$ !

Observe min-length superstring $S$ for $S_1, S_2$

is

$$S = \underbrace{AT}_{S_1}\underbrace{ATTT}_{S_2} \quad \overset{\wedge}{=} \quad \overset{2}{\underset{S_1 \longrightarrow S_2}{\bullet \longrightarrow}}$$

i.e. arc where
prefix-length min./
overlap max.

(2)   $S = \{$ ATAT , TATA, TATT , TTAT $\}$

$G_S$ :     P     (overlap : $\text{ov}(i,j) = |S_i| - p(i,j))$

$S_1 =$ ATAT•  ⟵ 40 ⟶ •  $S_2 =$ TATT
                    1 3

3│3              2² 2²      2│13
1│1              3ᵀ 3ᵀ      2│

$S_4 =$ TATA•  ⟵ 40 ⟶ •  $S_3 =$ TTAT
                    1 3

A   T   A   T   T   A   T   A

1:1 correspondence between orderings of
elements in $S' \subseteq S$ & paths in $G_S$.

EXMPL:

$S_1 =$ ATAT

             ↑
             1

$S_4 =$ TATA•  ⟵ 1 ⟵  •  $S_3 =$ TTAT

$S' = \{ S_1, S_3, S_4 \}$
ordering :   $S_3 \, S_4 \, S_1$

**DEF**   given order $S_{i_1} \ldots S_{i_\ell}$ (path in $G_3$) define:

merge $(S_{i_1} \ldots S_{i_\ell}) :=$

$$\text{pref}(i_1, i_2) \; \text{pref}(i_2, i_3) \ldots \text{pref}(i_{\ell-1}, i_\ell) \; S_{i_\ell}$$

$S_{i_\ell}$
$S_{i_{\ell-1}}$

$S_{i_2}$
$S_{i_1}$

merge $(S_{i_1} \ldots S_{i_\ell})$   $\text{pref}(i_1, i_2)$   $\text{pref}(i_{\ell-1}, i_\ell)$   $S_{i_\ell}$

**EXMPL (above)**   order $S_3, S_4, S_1$

$S_1 = \quad\quad\quad\quad$ T A T A
$S_4 = \quad\quad\quad$ A T A T
$S_3 = \quad$ T T A T
merge $(S_3 S_4 S_1):$ T T A T A T A

$\text{pref}(3,2)$
$\quad\quad\text{pref}(4,1)$
$\quad\quad\quad S_1$

By def: $|\text{merge}(S_{i_1}..S_{i_e})| = \sum\limits_{i=1}^{\ell-1} p(i, i+1) + |S_{i_e}|$

Basic Idea:

If for $S_i, S_j$, overlap $\text{ov}(S_i, S_j)$ is „large"
then $S_i, S_j$ might also overlap in genom

$\Rightarrow$ successively „merge" strings with large overlap.

GREEDY-SCS ($\mathcal{S}$)

WHILE ($|\mathcal{S}| > 1$)

$s, t \leftarrow$ 2 strings in $\mathcal{S}$ with max $\text{ov}(s, t)$

// s=t poss.ble

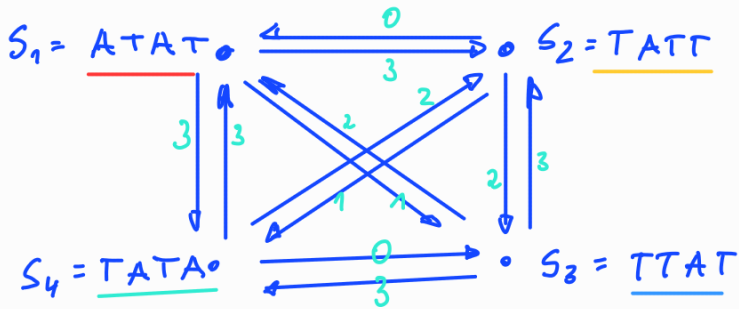remove $s, t$ from $\mathcal{S}$ & add $uvw$ to $\mathcal{S}$

// s= uv, t=vw

3. return the single string $s \in \mathcal{S}$.

EXMPL:

$\mathcal{S} = \{$ ATAT , TATA, TATT , TTAT $\}$

$G_\mathcal{S}:$ ov $(i,j)$



$S_1 = $ ATAT $\circ$ ⟵ 0 / 3 ⟶ $\circ$ $S_2 = $ TATT

$S_4 = $ TATA $\circ$ ⟶ 0 / 3 ⟵ $\circ$ $S_3 = $ TTAT
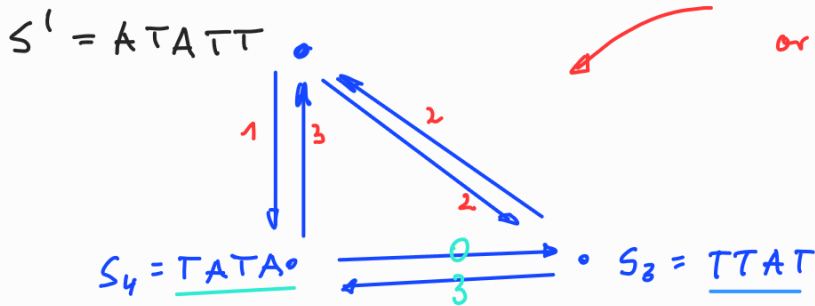
max overlap is 3    eg $(S_1, S_4)$, $(S_4, S_1)$, $(S_2, S_4)$, $(S_3, S_2)$...
take one of them & merge.
choose here $(S_1, S_2)$ :     $S_1 \neq S_2$

$\Rightarrow$  $\mathcal{S} = \{$ $\underset{\substack{\text{merge } (S_1, S_2) \\ \text{"} S'}}{\text{ATATT}}$ , $\underset{S_3}{\text{TTAT}}$ , $\underset{S_4}{\text{TATA}}$ $\}$

$S' = $ ATATT $\bullet$

ov. $(S', S_i) = $ ov $(S_2, S_i)$
ov $(S_i, S') = $ ov $(S_i, S_1)$



$S_4 = $ TATA $\circ$ ⟶ 0 / 3 ⟵ $\bullet$ $S_3 = $ TTAT

max ovulap: $(S_4, S')$ $\rightarrow$ merge $(S_4, S')$

get  $\mathcal{S} = ($ $\underset{S''}{\text{TATATT}}$ , $\underset{S_3}{\text{TTAT}}$ $)$

$S'' = $ TATATT $\bullet$

ov $(S'', S_3) = $ ov $(S_1, S_2)$
ov $(S_3, S'') = $ ov $(S_3, S_4)$

$S_3 = $ TTAT $\circ$

merge $(S_3, S'')$ :  TTATATT

<u>Lemma :</u>     GREEDY - SCS   has   runtime  $O(|S| \cdot \|\mathcal{J}\|)$

$$\|\mathcal{J}\| = \sum_{s \in \mathcal{S}} |s|$$

( Exercise , see also:

[Lit: Algorithmic Aspects of Bioinformatics, Böckenhauer & Bongartz 2007, springer ]  )

! important for the latter result is the following
result , that shows that overlaps in the
intermediate steps don't need to be recomputed !

<u>Lemma</u>    IF   in a step of Greedy - SCS . we merge $S'$ & $S''$

$$\text{\&} \quad S' = merge(S_{i_1} .. S_{i_k}) \quad k \geq 1$$
$$S'' = merge(S_{j_1} .. S_{j_\ell}) \quad \ell \geq 1$$

(prior to this step)

THEN

$$ov(S', S'') = ov(i_k, j_1)$$

<u>Proof:</u>   we show first that after merging $S'$ & $S''$ remains
"substring-free" ie   $\tilde{S} \subseteq S$ will never occur.

By assumption , before any merging,  $\mathcal{J}$ is substring-free

Assume, for contradiction, that after some merging
we have  $\tilde{S} \subseteq S$.

Let  $S$   be the first element constructed st
for some $\tilde{S} \in \mathcal{S}$ :  $\tilde{S} \subseteq S$

$S$ is either (a) the original $S$ (not part of any merging process so-far)

or (b) the result of merging $S_1, S_2$

(*$S_1, S_2$ could also have been merged before*)

[ Note $S \subseteq \tilde{S}$ not possible in case (a) since then $\exists$ not substring-free

& in case (b): $S_1, S_2 \subseteq \tilde{S}$ & $S$ not first element ] resp. $S$ not first element

it looks as follows:

$\Rightarrow \qquad S = \text{merge}(S_1, S_2)$



must look like

Since $\tilde{S} \subseteq S$ &

Since $S$ "first" $\Rightarrow \tilde{S} \nsubseteq S_1, S_2$

But then $\quad ov(S_1, \tilde{S}) \geq$
$ov(\tilde{S}, S_2) \geq \quad ov(S_1, S_2)$

& at least for one "$\geq$" we have "$>$"
& so greedy would choose $\tilde{S}$ & $S_i$, $i \in \{1, 2\}$
$\qquad \qquad \qquad$ ↯ greedy choice.

$\Rightarrow$ after each mergestep, $J$ remains substring free.

Assume now, for contradiction, that after some step $N \geq 1$, $S$ does not satisfy condition in the Lemma.

Put $S_{N-1} = S$ after applying $N-1$ steps of greedy.

By assumption on $S_{N-1}$ everything ok!

W.log let $S', S'' \in S_{N-1}$ sd $ov(S', S'') = $ max overlap over all elements in $S$

& greedy chooses $S'$ & $S''$ to be merged.

$\Rightarrow \quad S_N = S_{N-1} \setminus \{S', S''\} \cup \{\underbrace{merge(S', S'')}_{=S}\}$

Note, $S_N$ violates condition in lemma & in $S_N \setminus \{S\} \subseteq S_{N-1}$ everything is ok

$\Longrightarrow \quad \exists \; \tilde{S} \in S_N \setminus S$ st $S$ & $\tilde{S}$ violate condition in lemma.
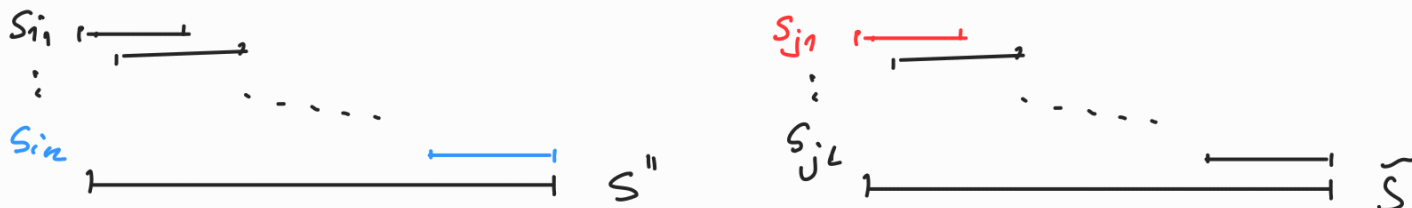
Situation is as follows:

$\qquad S = merge(S', S'')$

$\qquad S'' = merge(S_{i_1} .. S_{i_k}) \quad k \geq 1, \; S_{i_r} \in S_o$

$\qquad \tilde{S} = merge(S_{j_1} .. S_{j_\ell}) \quad \ell \geq 1, \; S_{j_s} \in S_o$

& since condition in lemma violated we have
$\qquad ov(S, \tilde{S}) \neq ov(S_{i_k}, S_{j_1})$
$[$ analog arguments for case $ov(\tilde{S}, S) \neq \cdots ]$

$S_{i_1}$, ... $S_{i_n}$ (blue), $S''$

$S_{j_1}$ (red), ... $S_{j_L}$, $\tilde{S}$

By definition: $S_{i_k} \sqsubseteq S = \text{merge}(S', S'')$

& $S_{j_1} \sqsubseteq \tilde{S}$

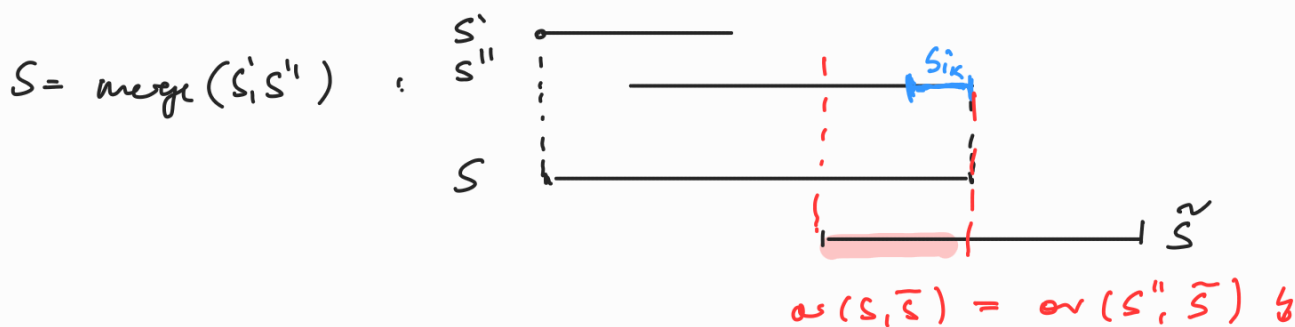$\Rightarrow \quad ov(S, \tilde{S}) \geq ov(S_{i_k}, S_{j_1})$

& since $" \neq "$

$\Rightarrow \quad ov(S, \tilde{S}) > ov(S_{i_k}, S_{j_1})$

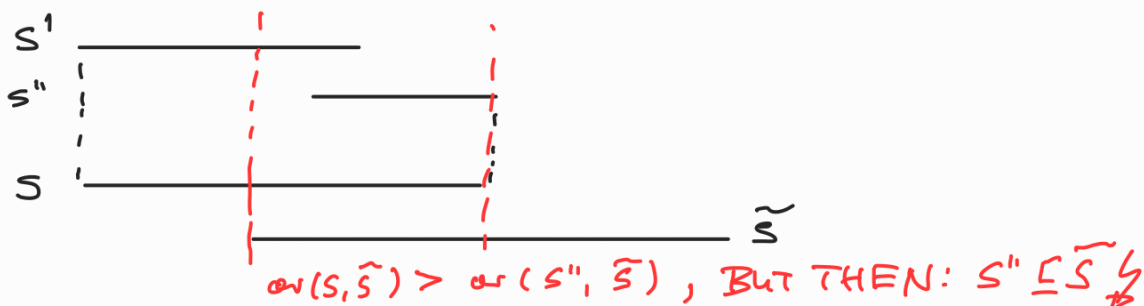By Ind-assumption (all true in $S_{N-1}$ & $\tilde{S}, S'' \in S_{N-1}$)
we have:

$$ov(S'', \tilde{S}) = ov(S_{i_k}, S_{j_1}) < ov(S, \tilde{S})$$

looks like:

Not Possible!

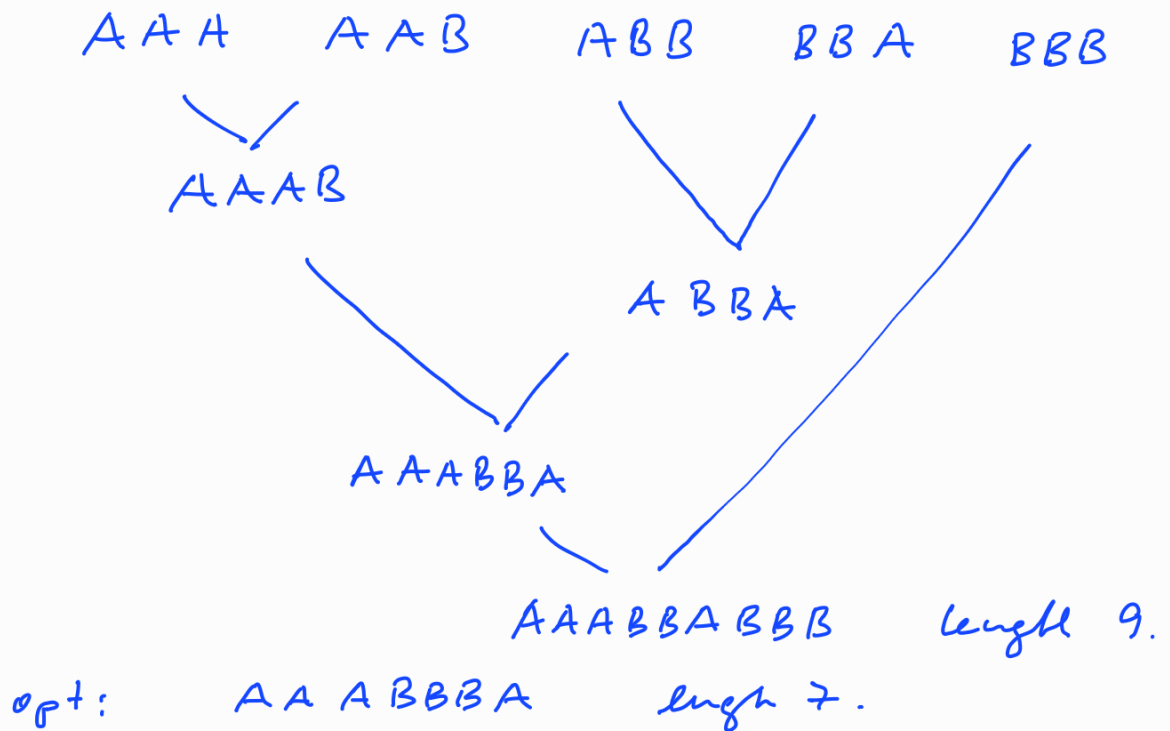$S = \text{merge}(S', S'')$

$S'$
$S''$

$S$

$S_{i_k}$

$\tilde{S}$

$ov(S, \tilde{S}) = ov(S'', \tilde{S})$ ⚡

$\Rightarrow$ must look like:

$S'$

$S''$

$S$

$\tilde{S}$

$ov(S, \tilde{S}) > ov(S'', \tilde{S})$, BUT THEN: $S'' \sqsubseteq \tilde{S}$ ⚡

In general, Greedy-SCS not optimal:

$$AAA \qquad AAB \qquad ABB \qquad BBA \qquad BBB$$

$$AAAB$$

$$ABBA$$

$$AAABBA$$

$$AAABBABBB \qquad \text{length } 9.$$

opt: $\qquad AAABBBA \qquad$ length 7.

But we are never "too far" away from opt. solution:

**Theorem:** String $S$ returned by Greedy-SCS satisfies:

$$|S| \leq 4\, opt(S)$$

ie, never worse than 4 times opt-solution.

[without proof]

One particular problem in practice: repeats.

In genomes often repeated regions.

Genomes often consist of repeated regions!

**Example:** Here, $\zeta$ = set of all substrings of size 6.

Greedy SCS on 6-mers of **a_long_long_long_time**

```
ng_lon _long_ a_long long_l ong_ti ong_lo long_t g_long g_time ng_tim
ng_time ng_lon _long_ a_long long_l ong_ti ong_lo long_t g_long
ng_time g_long_ ng_lon a_long long_l ong_ti ong_lo long_t
ng_time long_ti g_long_ ng_lon a_long long_l ong_lo
ng_time ong_lon long_ti g_long_ a_long long_l
ong_lon long_time g_long_ a_long long_l
long_lon long_time g_long_ a_long
long_lon g_long_time a_long
long_long_time a_long
a_long_long_time
```

The final superstring is shorter than the original "genome"

Genomes often consist of repeated regions!

**Example:** Here, $\zeta$ = set of all substrings of size 6.

& $|\zeta| = 16$ for all examples

```
a_long_long_time        a_long_long_long_time
a_long long_t           a_long long_l ng_tim
_long_ ong_ti           _long_ ong_lo g_time
long_l ng_tim           long_l ng_lon
long_l  g_time          ong_lo g_long
ong_lo                  ng_lon _long_
ong_lo                  g_long long_t
ng_lon                  _long_ ong_ti
ng_lon
g_long
g_long
_long_
_long_
```

```
a_long_long_long_long_long_time
a_long long_l  g_long    ng_tim
_long_   ng_lon long_l  g_time
ong_lo _long_   ng_lon
g_long ong_lo  _long_
long_t
ong_ti
```

To work with such problems one may employ: DeBruijn-graphs and Eulerian Paths. (board)

$k$-mer = substring of size $k$.

Given $J = \{$sequenced fragments$\}$

$\downarrow$

$J_k$ = all $k$-mers of sequences in $J$.

EXMPL:    (unknown) DNA    AAABBBBA

$\rightarrow$  $J = \{AAAA, AABB, BBA, ABBB\}$

$\rightarrow$  $J_3 = \{AAA, AAB, ABB, BBB, BBA\}$

DeBruijn-graph = multigraph (ie. multiple edges
$G_k$                    allowed)

Subdivide each $k$-mer $s_1 \ldots s_k$ into left $k-1$ mer $s_1 \ldots s_{k-1}$
right $k-1$-mer $s_2 \ldots s_k$

$\underbrace{\phantom{xxxxxxxxxxxxx}}$

these become
the vertices in
DeBruijn graph $G_k$

arc $(v, w)$ in $G_k$ if $v$ is L-$k-1$mer
$w$ is R-$k-1$mer
of some kmer
in $J_k$.

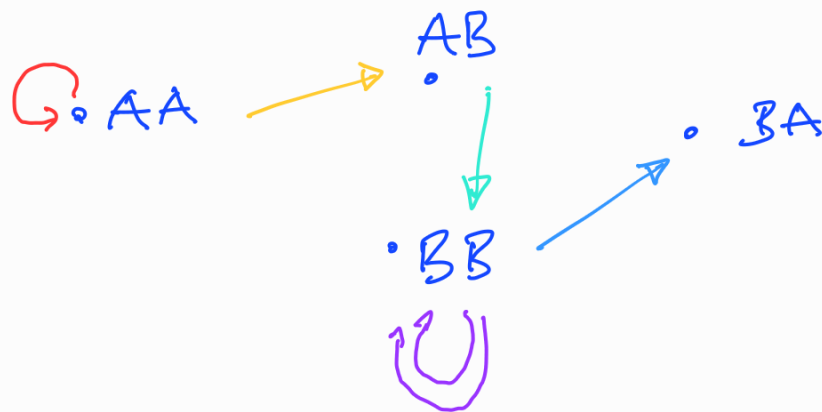Exmpl:     (unknown) DNA   AAABBBBA

all  3mers:   $S_3$ = { AAA, AAB, ABB, BBB, BBB, BBA }

can eg be obtained
from reads of fragments.
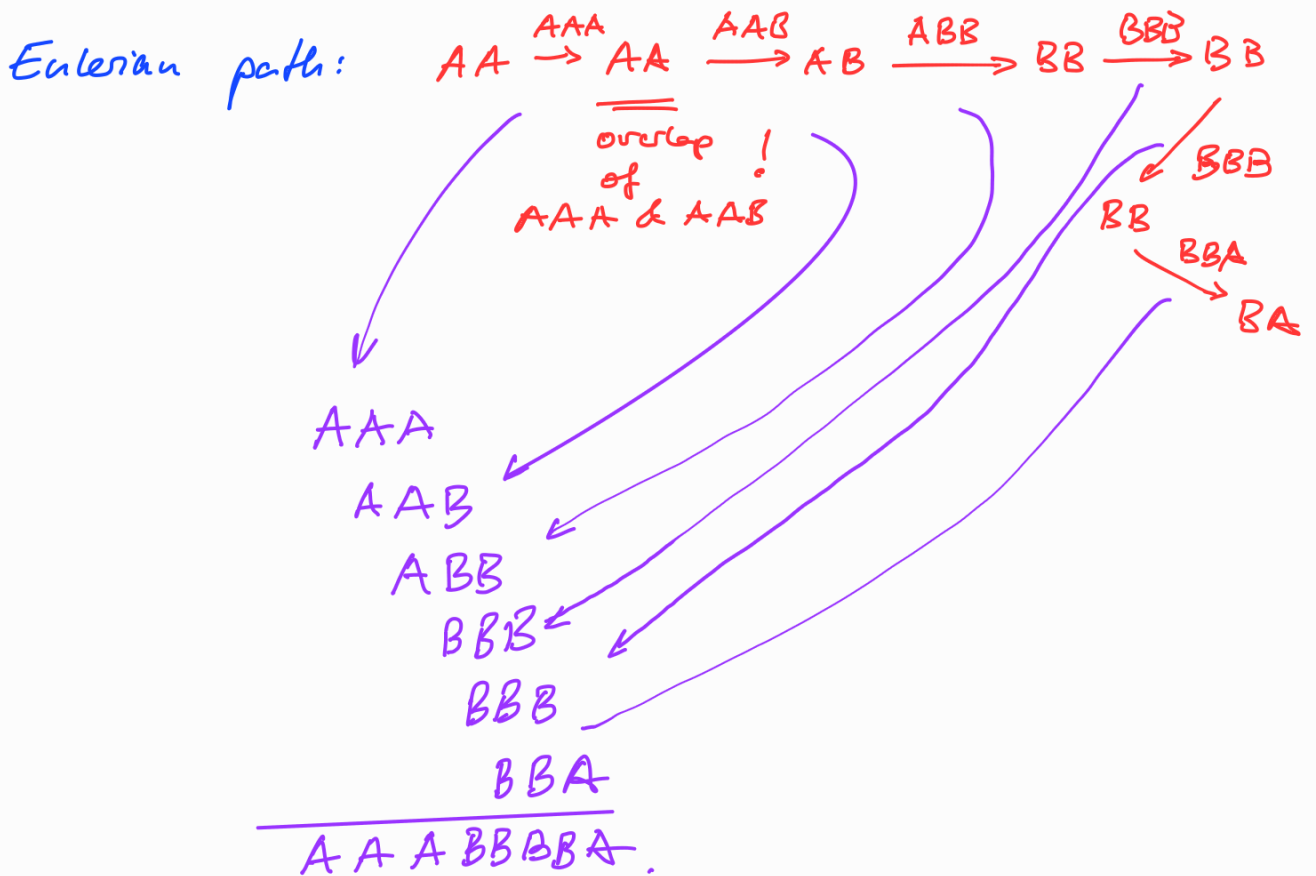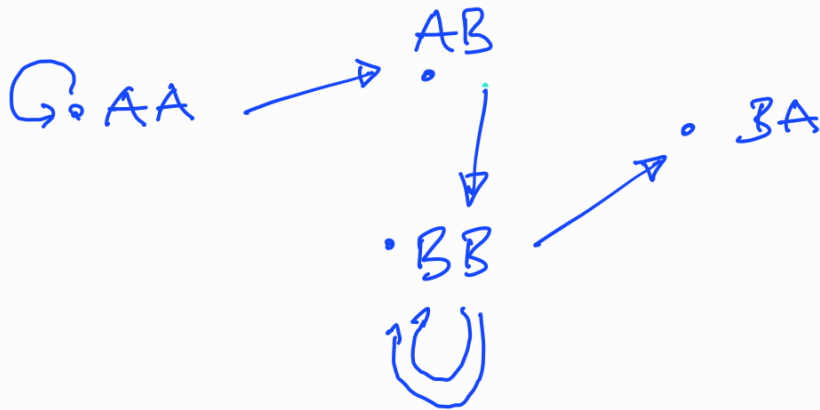
3-mers     AAA    →    L  k-1mer  AA
                       R  k-1mer  AA

           AAB         L   -"-    AA
                       R   -"-    AB

           ABB         L          AB
                       R          BB

      (2x) BBB         L          BB
                       R          BB

           BBA         L          BB
                       R          BA

Vertices in $S_3$ :

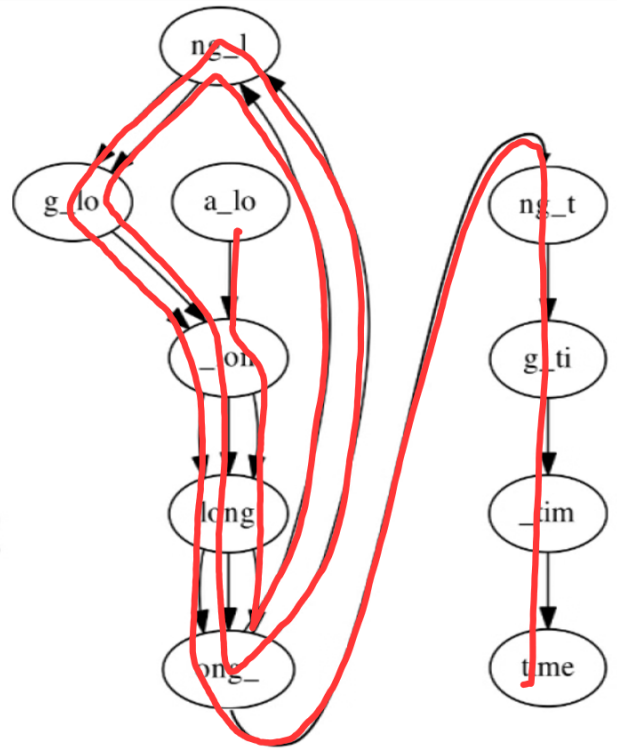

edges (implicitely) represent  3-mers   $AB \xrightarrow{ABB} BB$

Eulerian path in a directed graph is a path
that "visits" each edge exactly once.

G: AA $\rightarrow$ AB $\rightarrow$ BB $\rightarrow$ BA
     (AA loop, BB loop)

Eulerian path:

AA $\xrightarrow{AAA}$ AA $\xrightarrow{AAB}$ AB $\xrightarrow{ABB}$ BB $\xrightarrow{BBB}$ BB

overlap !
of
AAA & AAB

BB $\xrightarrow{BBA}$ BA

BBB

AAA
AAB
ABB
BBB
BBB
BBA
_____
A A A B B B B A .

De Bruijn graph (k=5) for:
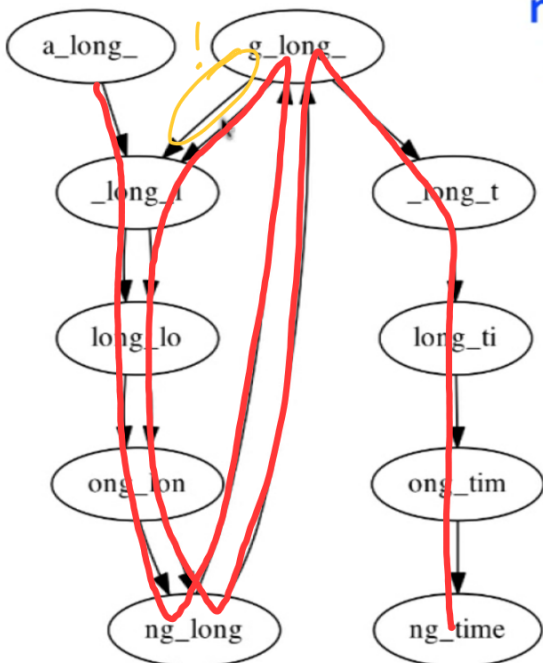
**a_long_long_long_time**

Eulerian walk gives original genome!



Eulerian Walk.

( typical k ~ 30-50 )

k = 8  Genome: a_long_long_long_time

Reads: a_long_long_long, ng_long_l, g_long_time

k-mers: a_long_l        ng_long_    g_long_t
        _long_lo        g_long_l    _long_ti
        long_lon                    long_tim
        ong_long                    ong_time
        ng_long_
        g_long_l
        _long_lo
        long_lon
        ong_long



No Eulerian walk!

Overlap graphs and DeBruijn graphs can be used to represent "relationships" between substrings.

The provided algorithms can, in general, not solve the assembly problem in an "optimal way" but serve as useful heuristics.

There are more sophisticated methods out there that are often based on these type pf algorithms.

---

[*]Medvedev & Pop *What do Eulerian and Hamiltonian cycles have to do with genome assembly?* PLoS Comput Biol. 2021

## Classical problems in practice:

► sequencing errors

► overlapping regions of fragments that are located on "far away" positions on DNA

► incomplete data (DNA not covered by resulting sequenced fragments)

► orientation of reads usually unknown

► repeats