

hhu.

Correcting Long-Reads with *k*-mers: A Dream Comes True

Pierre Marijon, Philipp Spohr, Antoine Limasset



November 23, 2020

- 1 Introduction
- 2 PanCov-Correct
- 3 Pcon & Br
- 4 Take home message

Long-read are usefull

Long-read are usefull

But they have a high error rate

Long-read are usefull

But they have a high error rate

Self-correction is efficient (low error rate) but:

Long-read are usefull

But they have a high error rate

Self-correction is efficient (low error rate) but:

- does not consider heterozygote variants

Long-read are usefull

But they have a high error rate

Self-correction is efficient (low error rate) but:

- does not consider heterozygote variants
- huge computation time (quadratic)

Long-read are usefull

But they have a high error rate

Self-correction is efficient (low error rate) but:

- does not consider heterozygote variants
- huge computation time (quadratic)

Remark: *k*-mers based methods work well on short-read and on hybrid correction data

- 1 Introduction
- 2 PanCov-Correct
- 3 Pcon & Br
- 4 Take home message

PanCov-Correct: Context

PanCov project goal: detect variants in COVID-19 samples

PanCov-Correct: Context

PanCov project goal: detect variants in COVID-19 samples

With unusual reads:

PanCov project goal: detect variants in COVID-19 samples

With unusual reads:

- reverse transcript amplified COVID-19

PanCov project goal: detect variants in COVID-19 samples

With unusual reads:

- reverse transcript amplified COVID-19
- Nanopore R9

PanCov project goal: detect variants in COVID-19 samples

With unusual reads:

- reverse transcript amplified COVID-19
- Nanopore R9
- 300bp !!! (due to lab protocol)

PanCov project goal: detect variants in COVID-19 samples

With unusual reads:

- reverse transcript amplified COVID-19
- Nanopore R9
- 300bp !!! (due to lab protocol)
- 300x coverage (coverage drop $\approx 20x$)

PanCov project goal: detect variants in COVID-19 samples

With unusual reads:

- reverse transcript amplified COVID-19
- Nanopore R9
- 300bp !!! (due to lab protocol)
- 300x coverage (coverage drop $\approx 20x$)
- $\approx 7\%$ error

PanCov project goal: detect variants in COVID-19 samples

With unusual reads:

- reverse transcript amplified COVID-19
- Nanopore R9
- 300bp !!! (due to lab protocol)
- 300x coverage (coverage drop $\approx 20x$)
- $\approx 7\%$ error
- strand bias

PanCov project goal: detect variants in COVID-19 samples

With unusual reads:

- reverse transcript amplified COVID-19
- Nanopore R9
- 300bp !!! (due to lab protocol)
- 300x coverage (coverage drop $\approx 20x$)
- $\approx 7\%$ error
- strand bias
- strain mixture

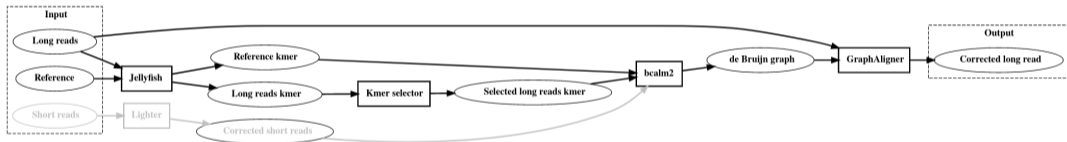
PanCov-Correct goal: correct reads while keeping variants, especially low-abundance strains

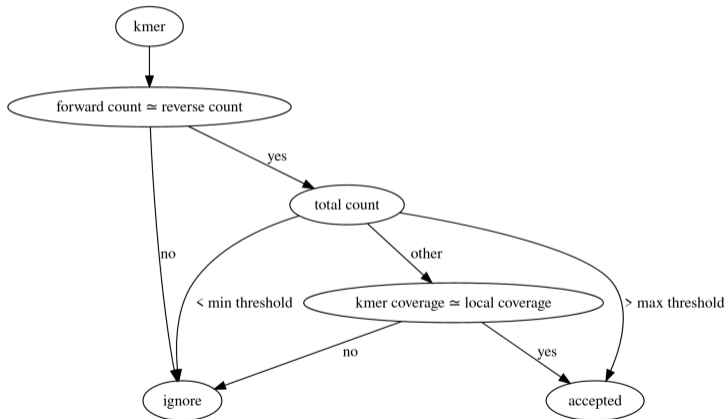
GraphAligner hybrid correction pipeline:



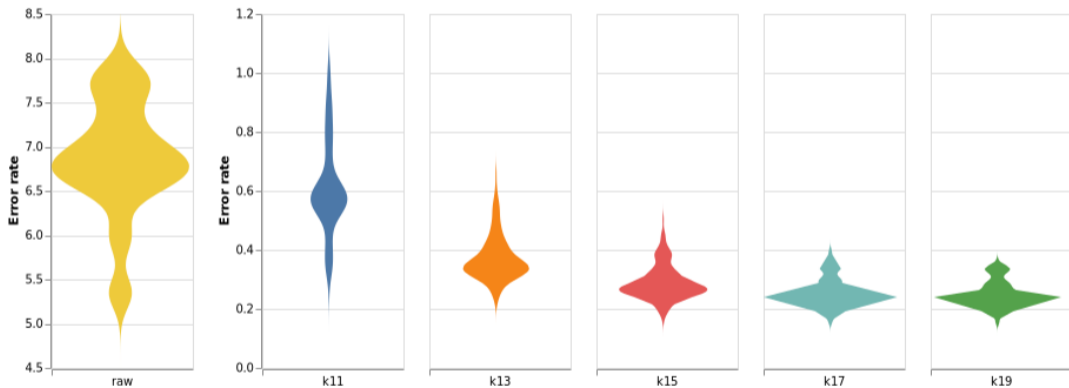
PanCov-Correct: Overview

GraphAligner hybrid correction pipeline:

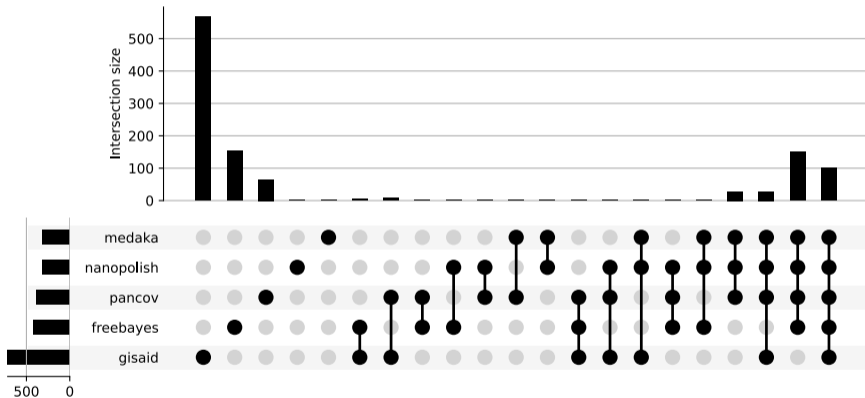


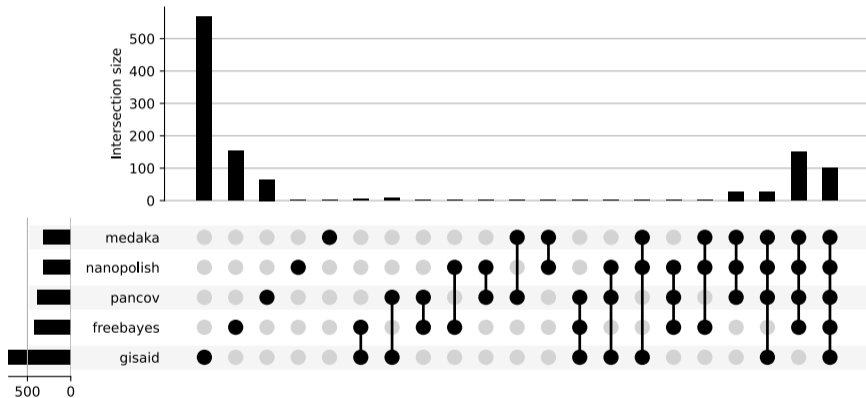


PanCov-Correct: Error rate



PanCov-Correct: Variant calling





Awaiting illumina sequencing for confirmation

On **our data** PanCov-Correct:

- Reduces error rate
- Retains low covered variants
- Retains *heterozygote* variants

On **our data** PanCov-Correct:

- Reduces error rate
- Retains low covered variants
- Retains *heterozygote* variants

Future:

- Standalone tools
- Reference free usage
- Test on other organisms and common reads
- Running time optimization

- 1 Introduction
- 2 PanCov-Correct
- 3 Pcon & Br**
- 4 Take home message

Pcon uses a hash function $k\text{-mer} \rightarrow [0, \frac{4^k}{2}[$, if k is odd

We define several functions:

- *kmer2bin* converts a DNA string into two bits $A \rightarrow 00$, $C \rightarrow 01$, $G \rightarrow 11$ and $T \rightarrow 10$
- *revcomp* performs the reverse complement for a binary representation of $k\text{-mer}$
- *popcount* count number of 1 in binary representation of a number

Pcon uses a hash function $k\text{-mer} \rightarrow [0, \frac{4^k}{2}[$, if k is odd

We define several functions:

- *kmer2bin* converts a DNA string into two bits $A \rightarrow 00$, $C \rightarrow 01$, $G \rightarrow 11$ and $T \rightarrow 10$
- *revcomp* performs the reverse complement for a binary representation of k -mer
- *popcount* count number of 1 in binary representation of a number

<code>kmer2bin(AGC)</code>	\rightarrow	<code>00 11 01</code>
<code>revcomp(00 11 01)</code>	\rightarrow	<code>11 01 10</code>
<code>popcount(00 11 01)</code>	\rightarrow	<code>3</code>
<code>popcount(11 01 10)</code>	\rightarrow	<code>4</code>

Pcon uses a hash function $k\text{-mer} \rightarrow [0, \frac{4^k}{2}[$, if k is odd

We define several functions:

- *kmer2bin* converts a DNA string into two bits A \rightarrow 00, C \rightarrow 01, G \rightarrow 11 and T \rightarrow 10
- *revcomp* performs the reverse complement for a binary representation of k -mer
- *popcount* count number of 1 in binary representation of a number

hash (*kmer*)

```
bin = kmer2bin(kmer)
if popcount(bin) % 2 == 0 then
  | return bin » 1
else
  | return revcomp(bin) » 1
```

```
kmer2bin(AGC)      → 00 11 01
revcomp(00 11 01) → 11 01 10
popcount(00 11 01) →          3
popcount(11 01 10) →          4
hash(AGC)          → 11 01 1
```

Pcon uses a hash function $k\text{-mer} \rightarrow [0, \frac{4^k}{2}[$, if k is odd

We define several functions:

- *kmer2bin* converts a DNA string into two bits $A \rightarrow 00$, $C \rightarrow 01$, $G \rightarrow 11$ and $T \rightarrow 10$
- *revcomp* performs the reverse complement for a binary representation of k -mer
- *popcount* count number of 1 in binary representation of a number

hash (*kmer*)

```
bin = kmer2bin(kmer)
if popcount(bin) % 2 == 0 then
  | return bin » 1
else
  | return revcomp(bin) » 1
```

count (*sequences*)

```
count = [0] ×  $\frac{4^k}{2}$ 
foreach sequence ∈ sequences do
  | foreach kmer ∈ sequence do
    | | count[hash(kmer)] += 1
return count
```

1 Count k -mers

- 1 Count k -mers
- 2 Analyse k -mers spectrum to find minimal abundance threshold

- 1 Count k -mers
- 2 Analyse k -mers spectrum to find minimal abundance threshold
- 3 Creation of a bitfield set.

- 1 Count k -mers
- 2 Analyse k -mers spectrum to find minimal abundance threshold
- 3 Creation of a bitfield set.
- 4 If k -mer count is higher than the minimal abundance threshold k -mer is added to bitfield.

- 1 Count k -mers
- 2 Analyse k -mers spectrum to find minimal abundance threshold
- 3 Creation of a bitfield set.
- 4 If k -mer count is higher than the minimal abundance threshold k -mer is added to bitfield.
- 5 Over each sequence apply correction algorithms

One: correct isolate error (musket algorithm with modification to support indel)

One: correct isolate error (musket algorithm with modification to support indel)

TGGTAGTAGTTACGA

Graph: search a simple path in *DeBruijn* graph between k -mer around error

One: correct isolate error (musket algorithm with modification to support indel)

TGGTAGTAGTTACGA

Graph: search a simple path in *DeBruijn* graph between k -mer around error

GapLength: use distance between k -mer around error to correct good number of base

One: correct isolate error (musket algorithm with modification to support indel)

TGGTAGTAGTTACGA

Graph: search a simple path in *DeBruijn* graph between k -mer around error

GapLength: use distance between k -mer around error to correct good number of base

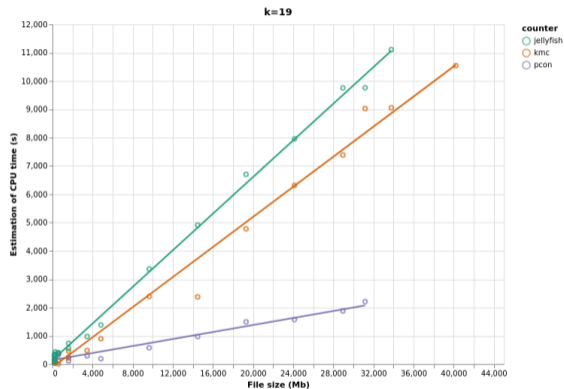
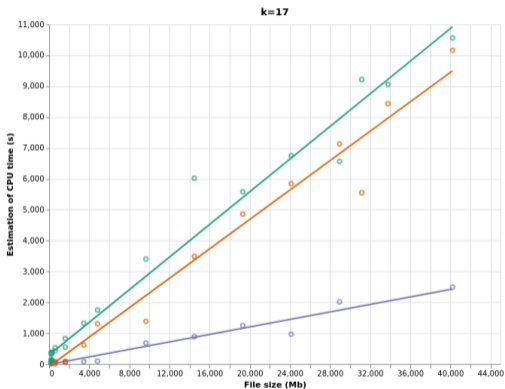
Greedy: get N in *DeBruijn* graph perform a pairwise alignment to check it's correct

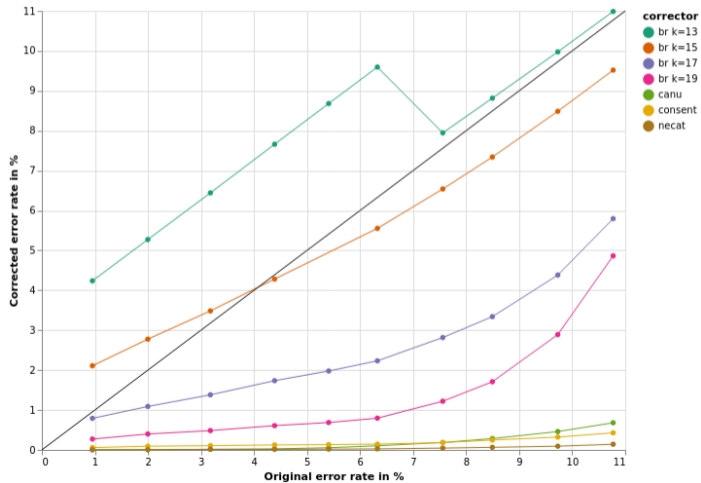
Codename	Organism	Technology	Error rate	Coverage
bacteria	<i>E. coli</i>	ONT R10.0	14.7%	≈ 127x
bacteria5	<i>E. coli</i>	ONT R10.3	5.9%	≈ 54x
bacteria7	<i>E. coli</i>	ONT R10.3	7.7%	≈ 127x
metagenome	metagenome	ONT R10.3	10.8%	
yeast	<i>S. cerevisiae</i>	ONT R10.3	8.3 %	≈ 283x
synthetic	<i>E. coli</i>	Badreads	*	≈ 50x
celegans	<i>C. elegans</i>	Badreads	5 %	**

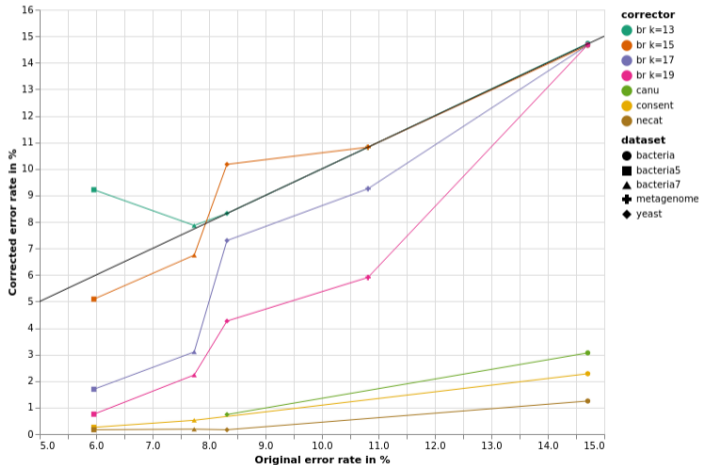
*: 1% to 10% per 1% step

** : 16x, 20x, 50x to 400x per 50x step

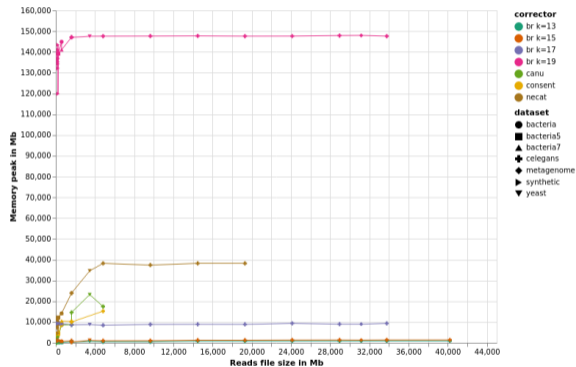
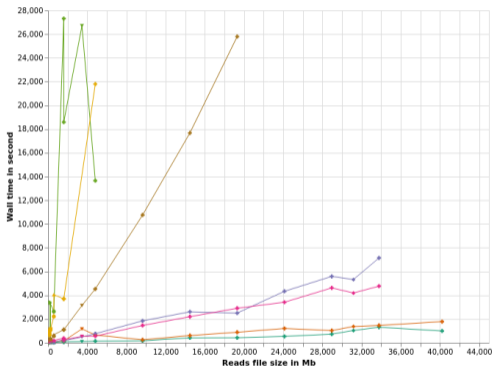
Pcon: Runtime and memory







Br: Runtime and memory



Conclusion: Pcon & Br

Pcon count k -mer faster than other tools but:

- only odd k -mer size
- small k -mer $k=13 \rightarrow 32\text{Mb}$ $k=21 \rightarrow 2\text{Tb}$
- only canonical form

Conclusion: Pcon & Br

Pcon count k -mer faster than other tools but:

- only odd k -mer size
- small k -mer $k=13 \rightarrow 32\text{Mb}$ $k=21 \rightarrow 2\text{Tb}$
- only canonical form

Br very fast but less efficient than other, future:

- other algorithms, **Two**, **Greedier**,
- use different set structure \rightarrow use larger kmer
- heterozygosity should be preserved, we have to check
- read filtering, scrubbing, contig polishing,
- hybrid correction

- 1 Introduction
- 2 PanCov-Correct
- 3 Pcon & Br
- 4 Take home message**

We can correct long-reads with long-reads k -mers:

- with long-read hybrid correction method \rightarrow PanCov-Correct
- with short-read correction method \rightarrow Pcon & Br
- do you have any new ideas?

We can correct long-reads with long-reads k -mers:

- with long-read hybrid correction method \rightarrow PanCov-Correct
- with short-read correction method \rightarrow Pcon & Br
- do you have any new ideas?

Improvement of raw reads quality will help us

We can correct long-reads with long-reads k -mers:

- with long-read hybrid correction method \rightarrow PanCov-Correct
- with short-read correction method \rightarrow Pcon & Br
- do you have any new ideas?

Improvement of raw reads quality will help us

No biorxiv, github link or bioconda logo here we are still in write/development

We can correct long-reads with long-reads k -mers:

- with long-read hybrid correction
- with short-read correction
- do you have any new ideas

Improvement of raw reads of

No biorxiv, github link or bio



development

Input: **TGGTAGTCCTTACGA**

✓ **TGGTA**

Input: **TGGTAGTCCTTACGA**

✓ **TGGTA**
✓ **GGTAG**

Input: **TGGTAGTCCTTACGA**

✓ **TGGTA**
✓ **GGTAG**
✗ **GTAGT**

Input: **TGGTAGTCCTTACGA**

- ✓ **TGGTA**
- ✓ **GGTAG**
- ✗ **GTAGT**
- ✗ **GTAGC**
- ✗ **GTAGA**
- ✓ **GTAGG**

Input: **TGGTAGTCCTTACGA**

- ✓ **TGGTA**
- ✓ **GGTAG**
- ✗ **GTAGT**
- ✗ **GTAGC**
- ✗ **GTAGA**
- ✓ **GTAGG**
- ✓ **TAGGC**
- ✓ **AGGCC**

Input: **TGGTAG****T****CCTTACGA**

- ✓ **TGGTA**
- ✓ **GGTAG**
- ✗ **GTAGT**
- ✗ **GTAGC**
- ✗ **GTAGA**
- ✓ **GTAGG**
- ✓ **TAGGC**
- ✓ **AGGCC**

Output: **TGGTAG****G****CCTTACGA**

Input: TGGTAG**T**CCTTACGA

✓ TGGTA
✓ GGTAG
✗ GTAGT
✗ GTAGC
✗ GTAGA
✓ GTAGG
✓ TAGGC
✓ AGGCC

Output: TGGTAG**G**CCTTACGAInput: TGGTAG**T****C**GTTACGA

✓ TGGTA
✓ GGTAG
✗ GTAGT
✗ GTAGC
✗ GTAGA
✓ GTAGG

Input: TGGTAG**T**CCTTACGA

✓ TGGTA
✓ GGTAG
✗ GTAGT
✗ GTAGC
✗ GTAGA
✓ GTAGG
✓ TAGGC
✓ AGGCC

Output: TGGTAG**G**CCTTACGAInput: TGGTAG**T****C**GTTACGA

✓ TGGTA
✓ GGTAG
✗ GTAGT
✗ GTAGC
✗ GTAGA
✓ GTAGG
✓ TAGGC
✗ AGGCC

Input: TGGTAG**T**CCTTACGA

✓ TGGTA
✓ GGTAG
✗ GTAGT
✗ GTAGC
✗ GTAGA
✓ GTAGG
✓ TAGGC
✓ AGGCC

Output: TGGTAG**G**CCTTACGAInput: TGGTAG**T****C**GTTACGA

✓ TGGTA
✓ GGTAG
✗ GTAGT
✗ GTAGC
✗ GTAGA
✓ GTAGG
✓ TAGGC
✗ AGGCC

Output: TGGTAG**T****C**GTTACGA

Input: TGGTAG**T**CCTTACGA

✓ TGGTA
 ✓ GGTAG
 ✗ GTAGT
 ✗ GTAGC
 ✗ GTAGA
 ✓ GTAGG
 ✓ TAGGC
 ✓ AGGCC

Output: TGGTAG**G**CCTTACGAInput: TGGTAG**T****C**GTTACGA

✓ TGGTA
 ✓ GGTAG
 ✗ GTAGT
 ✗ GTAGC
 ✗ GTAGA
 ✓ GTAGG
 ✓ TAGGC
 ✗ AGGCC

Output: TGGTAG**T****C**GTTACGAInput: TGGTAG**T**CCTTACGA

✓ TGGTA
 ✓ GGTAG
 ✗ GTAGT
 ✗ GTAGC
 ✓ GTAGA
 ✓ GTAGG

Input: TGGTAG**T**CCTTACGA

✓ TGGTA
 ✓ GGTAG
 ✗ GTAGT
 ✗ GTAGC
 ✗ GTAGA
 ✓ GTAGG
 ✓ TAGGC
 ✓ AGGCC

Output: TGGTAG**G**CCTTACGAInput: TGGTAG**T****C**GTTACGA

✓ TGGTA
 ✓ GGTAG
 ✗ GTAGT
 ✗ GTAGC
 ✗ GTAGA
 ✓ GTAGG
 ✓ TAGGC
 ✗ AGGCC

Output: TGGTAG**T****C**GTTACGAInput: TGGTAG**T**CCTTACGA

✓ TGGTA
 ✓ GGTAG
 ✗ GTAGT
 ✗ GTAGC
 ✓ GTAGA
 ✓ GTAGG

Output: TGGTAG**T**CCTTACGA

Input: TGGTAG**T**CCTTACGA

✓ TGGTA
 ✓ GGTAG
 ✗ GTAGT
 ✗ GTAGC
 ✗ GTAGA
 ✓ GTAGG
 ✓ TAGGC
 ✓ AGGCC

Output: TGGTAG**G**CCTTACGAInput: TGGTAG**T****C**GTTACGA

✓ TGGTA
 ✓ GGTAG
 ✗ GTAGT
 ✗ GTAGC
 ✗ GTAGA
 ✓ GTAGG
 ✓ TAGGC
 ✗ AGGCC

Output: TGGTAG**T****C**GTTACGAInput: TGGTAG**T**CCTTACGA

✓ TGGTA
 ✓ GGTAG
 ✗ GTAGT
 ✗ GTAGC
 ✓ GTAGA
 ✓ GTAGG

Output: TGGTAG**T**CCTTACGA

$$\mathcal{O}(\text{set access}) = 4 + N$$

Input: **TGGTAGTAGTTACGA**

GGTAG **TTACG**

Input: **TGGTAGTAGTTACGA**

✓ **GGTAG** **TTACG**
✓ **GTAGG**
✓ **TAGGA**
✓ **AGGAC**
✓ **GGACT**
✓ **GACTT**
✓ **ACTTA**
✓ **CTTAC**
✓ **TTACG**

Output: **TGGTAGGACTTACGA**

Input: **TGGTAGTAGTTACGA**

✓ **GGTAG** **TTACG**
✓ **GTAGG**
✓ **TAGGA**
✓ **AGGAC**
✓ **GGACT**
✓ **GACTT**
✓ **ACTTA**
✓ **CTTAC**
✓ **TTACG**

Output: **TGGTAGGACTTACGA**

Criteria to stop graph exploration:

- number of successor $\neq 1$
- back on a k -mer seen before

Input: **TGGTAGTAGTTACGA**

✓
✓
✓
✓
✓
✓
✓
✓

GGTAG **TTACG**
GTAGG
TAGGA
AGGAC
GGACT
GACTT
ACTTA
CTTAC
TTACG

Output: **TGGTAGGACTTACGA**

Criteria to stop graph exploration:

- number of successor $\neq 1$
- back on a k -mer seen before

$$|\text{set access}| = 2 \times 8 = 2 \times (5 + 3)$$

Input: **TGGTAGTAGTTACGA**

✓
✓
✓
✓
✓
✓
✓
✓

GGTAG **TTACG**
GTAGG
TAGGA
AGGAC
GGACT
GACTT
ACTTA
CTTAC
TTACG

Output: **TGGTAGGACTTACGA**

Criteria to stop graph exploration:

- number of successor $\neq 1$
- back on a k -mer seen before

$$|\text{set access}| = 2 \times 8 = 2 \times (5 + 3)$$

$$\mathcal{O}(\text{set access}) = 2 \times (k\text{-mer size} + \text{error length})$$

Br: Method `GapLength`

Input: **TGGTAGTAGTTACGA**

GGTAG **TTACG**

Input: **TGGTAGTAGTTACGA**
GGTAG **TTACG**

```
GapLength(begin, distance)
if distance == kmer size then
  | One()
else if distance < kmer size then
  | Graph()
else
  | get_kmers(begin, distance - kmer size)
```

If distance > k-mer size:
 $\mathcal{O}(\text{set access}) = k\text{-mer size} + 2 \times \text{error length}$

Input: **TGGTAGTAGTTACGA**

GGTAG **TTACG**

✓
✓
✓
GTAGG

TAGGA

AGGAC

GapLength(*begin*, *distance*)

if *distance* == *kmer size* **then**

 | One()

else if *distance* < *kmer size* **then**

 | Graph()

else

 | get_kmers(*begin*, *distance* - *kmer size*)

If *distance* > *k-mer size*:

$\mathcal{O}(\text{set access}) = k\text{-mer size} + 2 \times \text{error length}$

Input: TGGTAGTAGTTACGA

GGTAG TTACG
✓ GTAGG
✓ TAGGA
✓ AGGAC

Output: TGGTAGGACTTACGA

```
GapLength(begin, distance)
if distance == kmer size then
  | One()
else if distance < kmer size then
  | Graph()
else
  | get_kmers(begin, distance - kmer size)

If distance > k-mer size:
 $\mathcal{O}(\text{set access}) = k\text{-mer size} + 2 \times \text{error length}$ 
```

Br: Method Greedy

Input: **TGGTAGTAGTTACGA**

Input: **TGGTAGTAGTTACGA**

✓ **GTAGG**
✓ **TAGGA**
✓ **AGGAC**
✓ **GGACT**
✓ **GACTT**

Input: TGGTAGTAGTTACGA

✓ GTAGG
✓ TAGGA
✓ AGGAC
✓ GGACT
✓ GACTT

GTAGTAGTT
| | | | | | |
GTAGGACTT

Input: **TGGTAGTAGTTACGA**

✓ **GTAGG**
✓ **TAGGA**
✓ **AGGAC**
✓ **GGACT**
✓ **GACTT**

GTAGTAGTT
| | | | | | | |
GTAGGACTT

Output: **TGGTAGGACTTACGA**

$$\mathcal{O}(\text{set access}) = M + \mathcal{O}(2 \times (K+M))$$