

May 27-29, 2026

Rectorate of the University of Roma Tre, Rome, Italy

STORAGE AND COMPUTING WITH **Dna**

BBQ: Improving Consensus Finding in DNA Data Storage with Base-Level Quality Scores

Puru Sharma

National University of Singapore

Sponsored by



European
Innovation
Council



Organized by



Introduction



The amount of data is growing exponentially, and we are running out of silicon to store it.

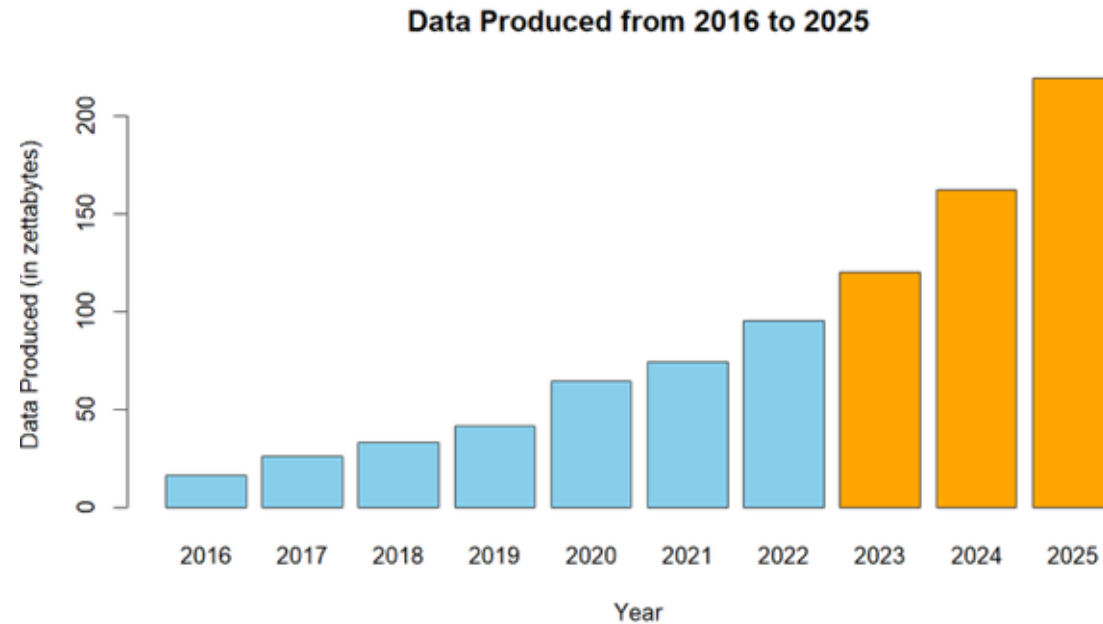


Figure. The amount of data produced each year [Hassini et al., 2023].

Zettabyte = 10^{21} bytes = 1,000,000,000,000,000,000 bytes



Introduction – Why DNA data storage?

1. Incredible density
~ 6 billion bits of info in human cells.
2. Unmatched durability
3. Little Energy Consumption
4. Constantly improving reading/writing interface.

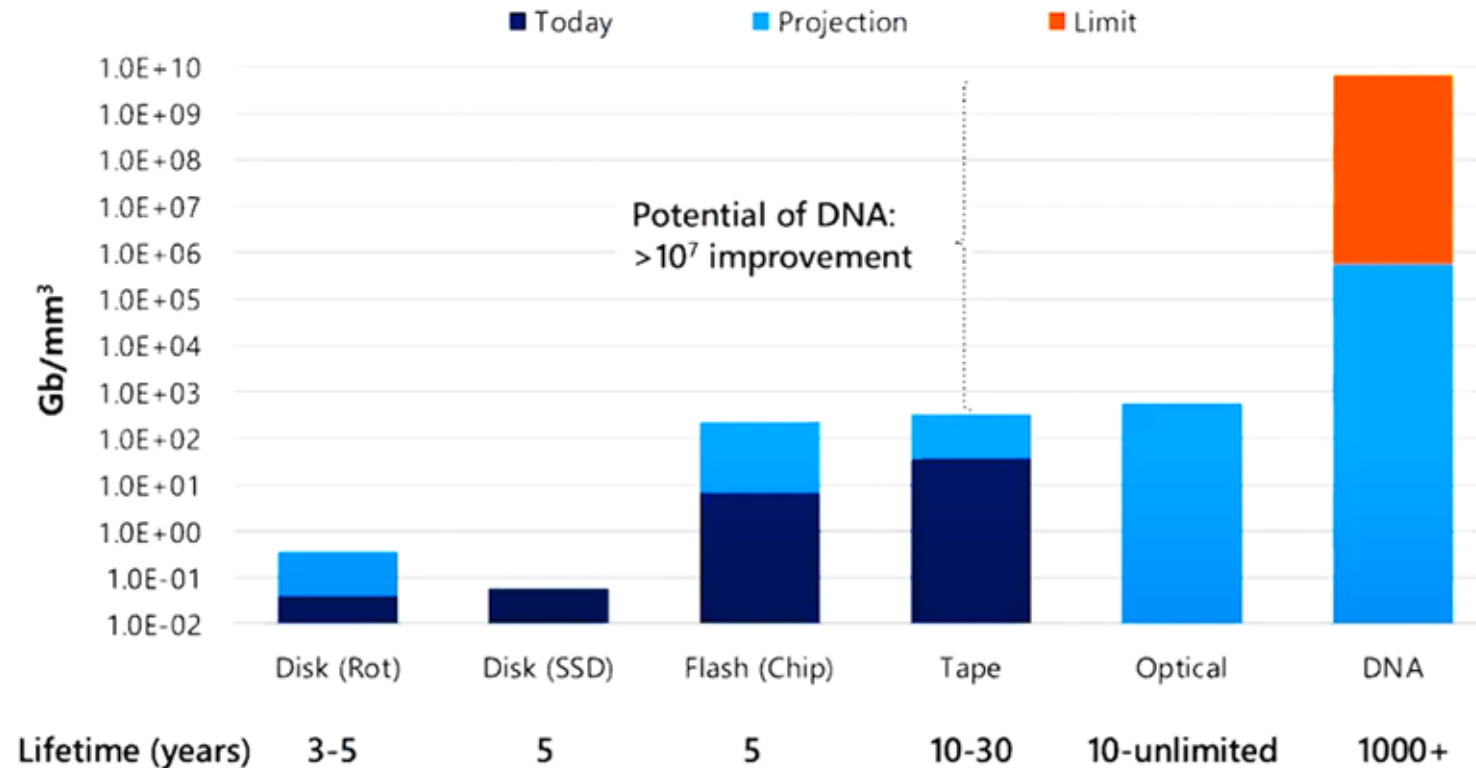


Figure. Storage density and lifetime.

Credit: Keynote by Luis Ceze and Karin Strauss, FAST '21.



Introduction – DNA Storage Pipeline: Writing



Digital data

010010001101
010010101000

Binary

00→A, 01→C,
10→G, 11→T.



CAGATC
CAGGGA

DNA Code

Introduction – DNA Storage Pipeline: Writing



Digital data

010010001101
010010101000

Binary

00→A, 01→C,
10→G, 11→T.



CAGATC
CAGGGA

DNA Code



AACAGA
ACTCCA
AGGGGA

Encoded strands

Typically in length 100-200.



Synthesized DNA

Introduction – DNA Storage Pipeline: Reading



AGTCCA AGGGA
AATAGA GGGGA ACTTCCA
ATCTAGC ACTCCCA
 AACGA AACGA

PCR and Sequencing

*With Illumina or
Oxford Nanopore*

Introduction – DNA Storage Pipeline: Reading



AGTCCA	AGGGA	AACGA
ACTTCCA	GGGGA	AATAGA
ACTCCCA	AGGCGA	AACGA

PCR and Sequencing

*With Illumina or
Oxford Nanopore*

Clustered Reads

*Reads that come from the same
template are clustered together.*

Introduction – DNA Storage Pipeline: Reading



PCR and Sequencing

*With Illumina or
Oxford Nanopore*



AGTCCA	AGGGA	AACGA
ACTTCCA	GGGGA	AATAGA
ACTCCA	AGGCGA	AACGA
↓	↓	↓
ACTCCA	AGGGGA	AACAGA

Trace Reconstruction

Consensus construction of each cluster.



Introduction – DNA Storage Pipeline: Reading



PCR and Sequencing

With Illumina or
Oxford Nanopore

AGTCCA AGGGA AACGA
ACTTCCA GGGGA AATAGA
ACTCCA AGGCGA AACGA

↓ ↓ ↓
ACTCCA AGGGGA AACAGA

CAGATCCAGGGA

← 010010001101010010101000

Decoding

Full decoding of
the original data



Introduction – Trace Reconstruction Problem

Problem: Given N traces, the noisy copies of original strand x of length L , we aim to reconstruct x with high probability.

$x = \text{GTAGTGCCTG}$

$y_1 = \text{GTAGGTGCCG}$

$y_2 = \text{GTAGTCCTG}$

$y_3 = \text{GTAGTGCCTG}$

$y_4 = \text{GTAGCGCCAG}$

$y_5 = \text{GCATGCTCTG}$



GTAGGTGCC-G

GTA-GT-CCTG

GTA-GTGCCTG

GTA-GCGCCAG

GCATGCT-CTG

$\hat{x} = \text{GTA-GTGCCTG}$

1. Cope with high error rate

- Errors can be introduced during synthesis, sequencing and clustering.

2. Deal with low coverage.

- Some clusters may be noisy and have very few traces.

3. Utilize prior knowledge.

- The target sequence length L is known.
- The sequencing technology provides us a quality score for each base called.



Method	Representative Tools
Alignment	<ul style="list-style-type: none">• MUSCLE [Edgar, 2004],• CPL [Bar-Lev et al., 2025].
Read error correction	<ul style="list-style-type: none">• BMA [Batu et al., 2004],• Trellis BMA [Srinivasavaradhan et al., 2021].
Assembly	<ul style="list-style-type: none">• DBGPS [Song et al., 2022],• ITR [Sabary et al., 2024].
Deep learning	<ul style="list-style-type: none">• DNA-GAN [Zheng et al., 2024],• DNAFormer [Bar-Lev et al., 2025].



Potential problems:

- Longer running time.
- Multiple ways of defining the alignment scoring scheme.

Assuming $p_{\text{insertion}} = p_{\text{deletion}} = p_{\text{substitution}}?$

AGTCCACT

AGAACTT

AGTCCCATT

Assuming $p_{\text{substitution}} = 0.001$,
 $p_{\text{insertion}} = p_{\text{deletion}} = 0.01?$



Potential problems:

- Longer running time.
- Multiple ways of defining the alignment scoring scheme.

AGTCCACT

$$\begin{cases} \text{match} = 1 \\ \text{mismatch} = -1 \\ \text{indel} = -1 \end{cases}$$

AGAACTT

AGTCCCATT

$$\begin{cases} \text{match} = 1 \\ \text{mismatch} = -1 \\ \text{indel} = -0.5 \end{cases}$$

Potential problems:

- Hard to capture complications in real data.

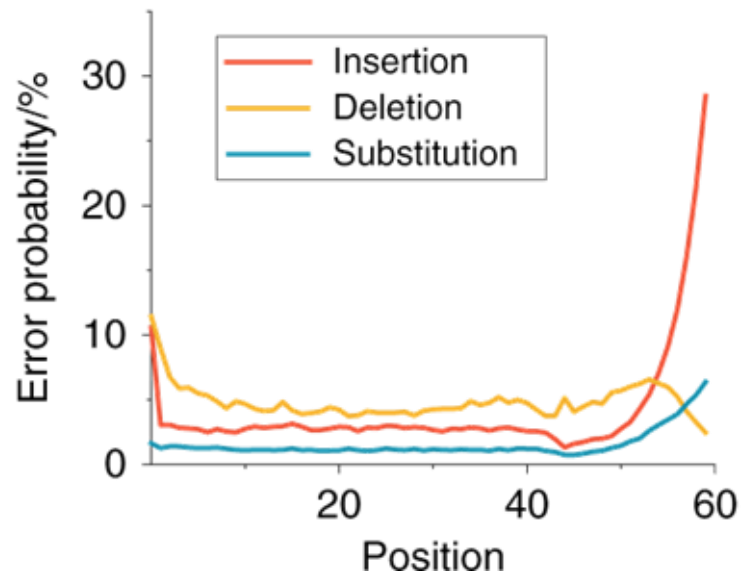


Figure. Error distribution within a read [Antkowiak et al., 2020].

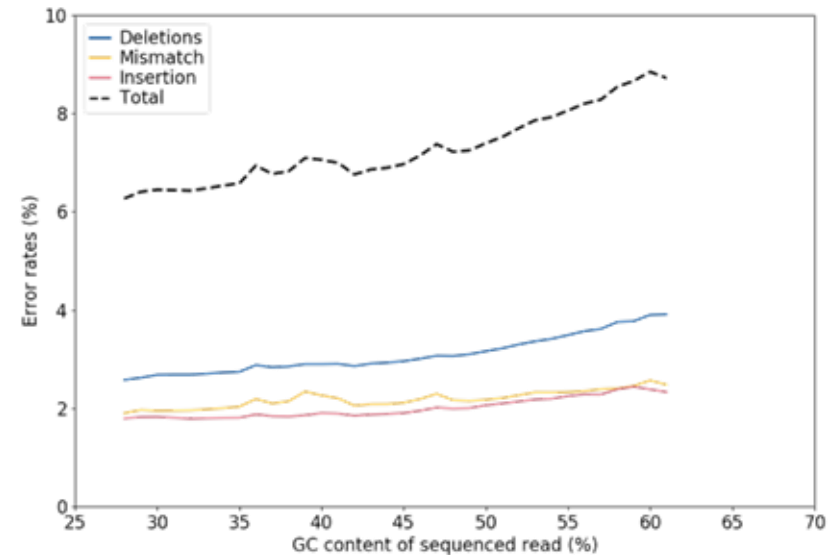


Figure. Higher GC-content may lead to higher error rates [Delahaye et al., 2021].

Potential problems:

- Hard to capture complications in real data.

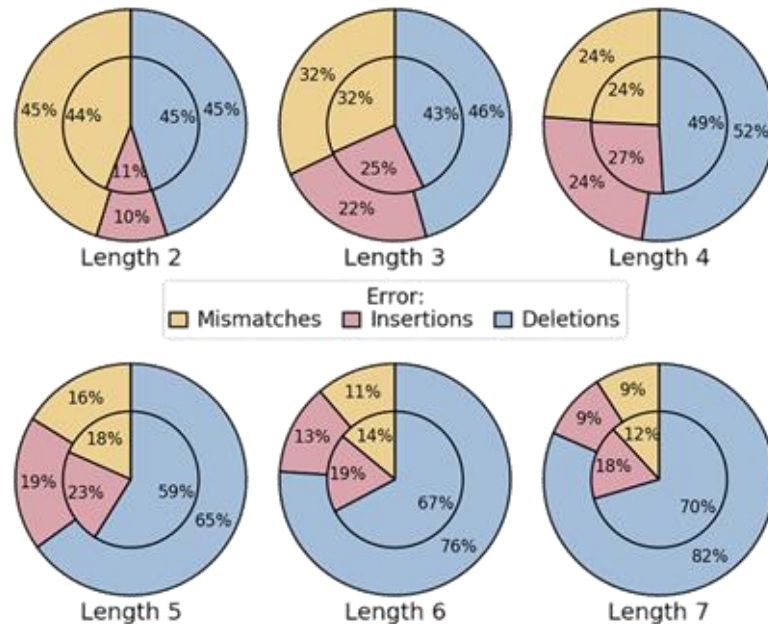


Figure. Error distribution in homopolymer regions [Delahaye et al., 2021].

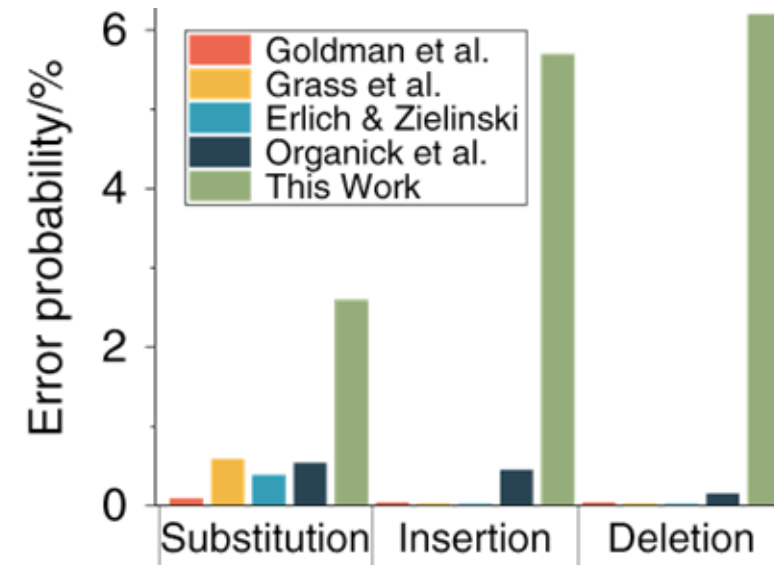


Figure. Error composition in various datasets [Antkowiak et al., 2020].

Potential problems:

- Majority aren't always correct, especially given the prior knowledge of the length of x .

AGTCCACT	$\left\{ \begin{array}{l} \text{match} = 1 \\ \text{mismatch} = -1 \\ \text{indel} = -1 \end{array} \right.$	AGTCC-ACT
AGAACTT		AGAAC--TT
AGTCCCATT	$\left\{ \begin{array}{l} \text{match} = 1 \\ \text{mismatch} = -1 \\ \text{indel} = -0.5 \end{array} \right.$	AGTCCCATT
		AGTCC-ATT
		AGTCC-ACT-
		AGA---ACTT
		AGTCCCA-TT
		AGTCC-ACTT

Potential problems:

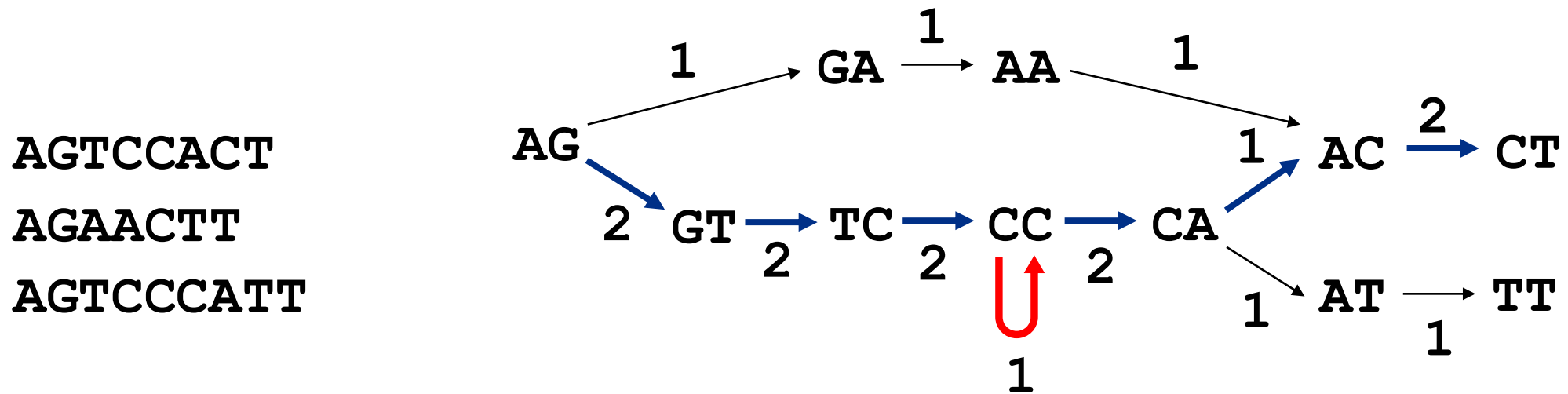
- Majority aren't always correct, especially given the prior knowledge of the length of x .

AGTCCACT	$\left\{ \begin{array}{l} \text{match} = 1 \\ \text{mismatch} = -1 \\ \text{indel} = -1 \end{array} \right.$	AGTCC-ACT	Given $L = 10$,	
AGAACTT		AGAAC--TT		AGTCC-ACT-
AGTCCCATT	$\left\{ \begin{array}{l} \text{match} = 1 \\ \text{mismatch} = -1 \\ \text{indel} = -0.5 \end{array} \right.$	AGTCCCATT		AGA---ACTT
		<u>AGTCC-ATT</u>		AGTCCCA-TT
		AGTCC-ACT-	<u>AGTCCCACTT</u>	
		AGA---ACTT		
		AGTCCCA-TT		
		<u>AGTCC-ACTT</u>		



Potential problems:

- Majority aren't always correct, especially given the prior knowledge of the length of x .



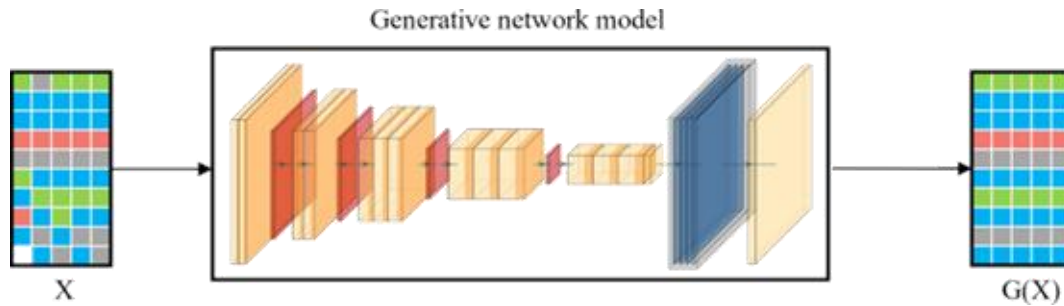


Figure. DNA-GAN [Zheng et al., 2024].

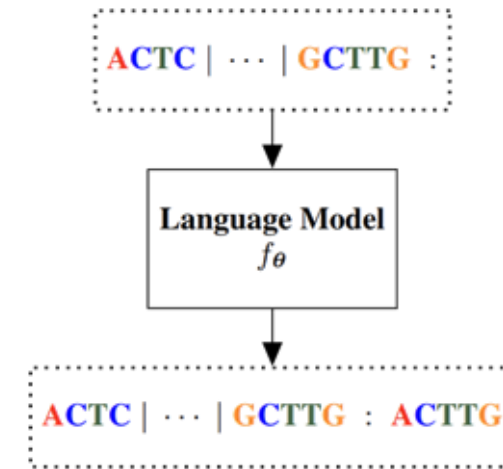


Figure. TRReconLM. [Girsch et al., 2025].

- **Problems:** Need for training; may not work equally well on unseen error distributions.

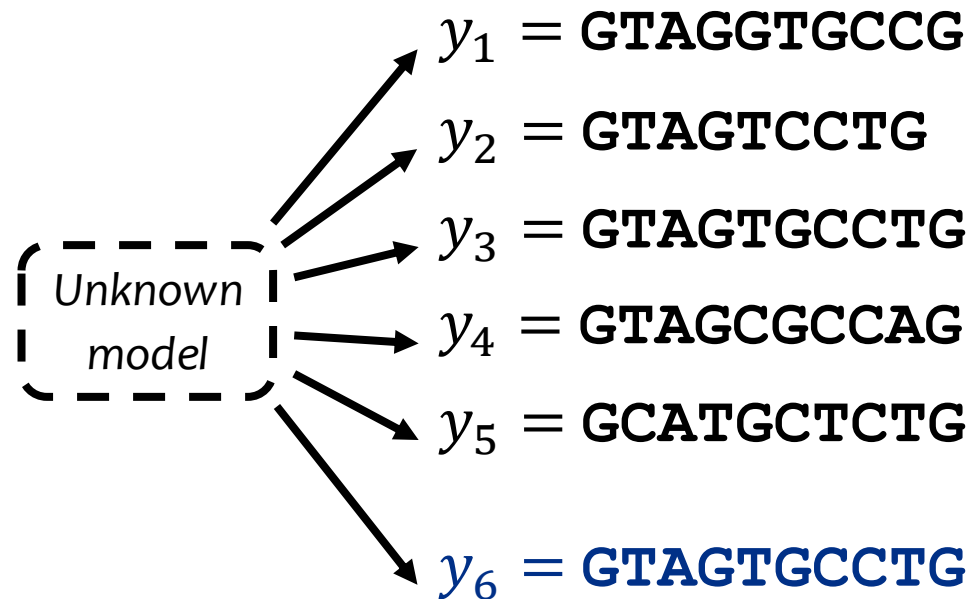
Method



GTAGGTGCC-G
GTA-GT-CCTG
GTA-GTGCCTG
GTA-GCGCCAG
GCATGCT-CTG

$\hat{x} = \text{GTA-GTGCCTG}$

What is the best alignment?



What is the most probable next trace?

- **Idea:** Model the traces as observations of a k -th order Markov chain.



Original Traces \mathcal{C}

```
AAAGTCG
AAAGTCGG
GAAGCG
GAAGGTC
GAAGTCG
GGAAGTTG
```

- k is chosen such that all the k -mers in the reads are unique.
- Estimate the parameters of the Markov chain using maximum likelihood.

$$\Pr[s_{1:4} = AAAG] = \frac{2}{6}$$

$$\Pr[s_{1:4} = GAAG] = \frac{3}{6}$$

Original Traces \mathcal{C}

```
AAAGTCG
AAAGTCGG
GAAGCG
GAAGGTC
GAAGTCG
GGAAGTTG
```



($k+1$)-mer profile

AAAG: 2	AGTC: 3
GAAG: 3	AGTT: 1
GGAA: 1	GTCG: 3
AAGT: 4	...

$$\begin{aligned} & \Pr[s_i = C \mid s_{i-k:i-1} = \text{AGT}] \\ &= \frac{n(\text{AGTC}) + 1}{n(\text{AGTA}) + n(\text{AGTC}) + n(\text{AGTG}) + n(\text{AGTT}) + 4} \\ &= \frac{4}{8} \end{aligned}$$



Original Traces \mathcal{C}

```
AAAGTCG
AAAGTCGG
GAAGCG
GAAGGTC
GAAGTCG
GGAAGTTG
```

$(k+1)$ -mer profile

AAAG: 2	AGTC: 3
GAAG: 3	AGTT: 1
GGAA: 1	GTCG: 3
AAGT: 4	...

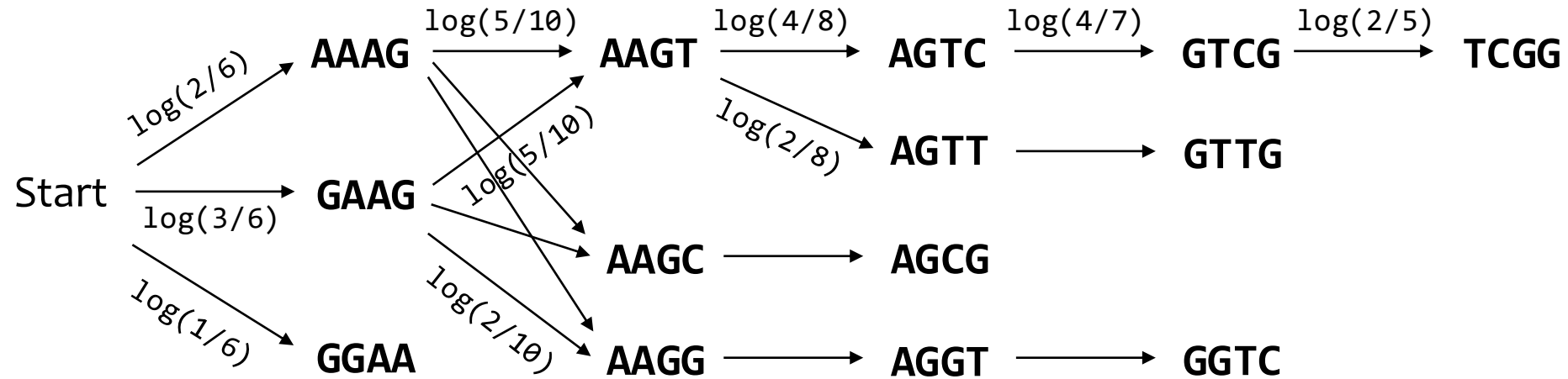
$$\Pr[s_i = C \mid s_{i-k:i-1} = AGT] = \frac{4}{8}$$
$$\Pr[s_i = T \mid s_{i-k:i-1} = AGT] = \frac{2}{8}$$

$\frac{2}{8}$ chance of substitution to T

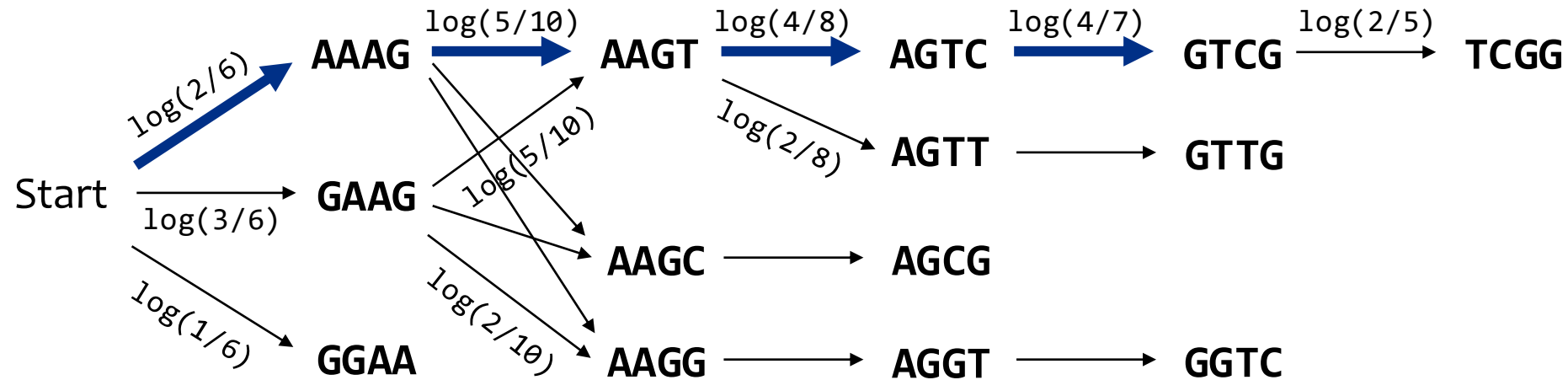
- Intuitively, the parameters of the Markov chain captures position-specific error profiles.



Substitute weights in the de Bruijn graph with the learned parameters of the Markov chain.



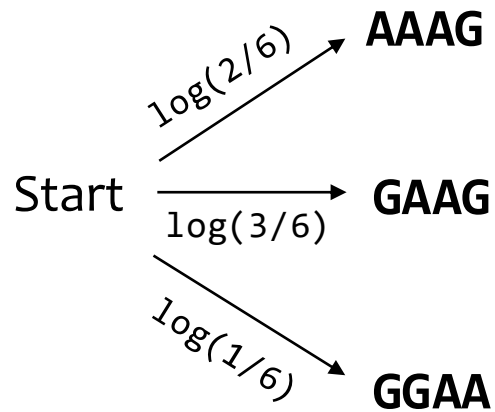
Goal: find the highest-weight path of length $L - k$ in the modified de Bruijn graph.



Sum of edge weight = $\log \Pr[S = \text{AAAGTCG}]$.



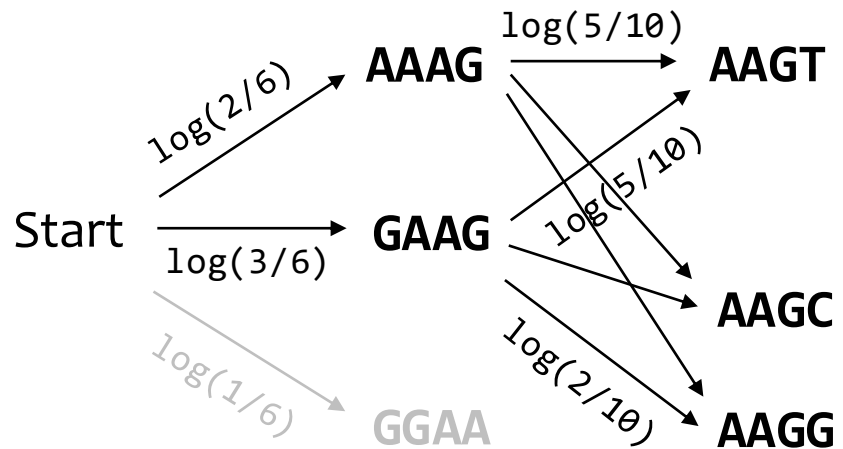
Goal: find the highest-weight path of length $L - k$ in the modified de Bruijn graph.



Beam search: keep only the top B paths (here $B = 2$).



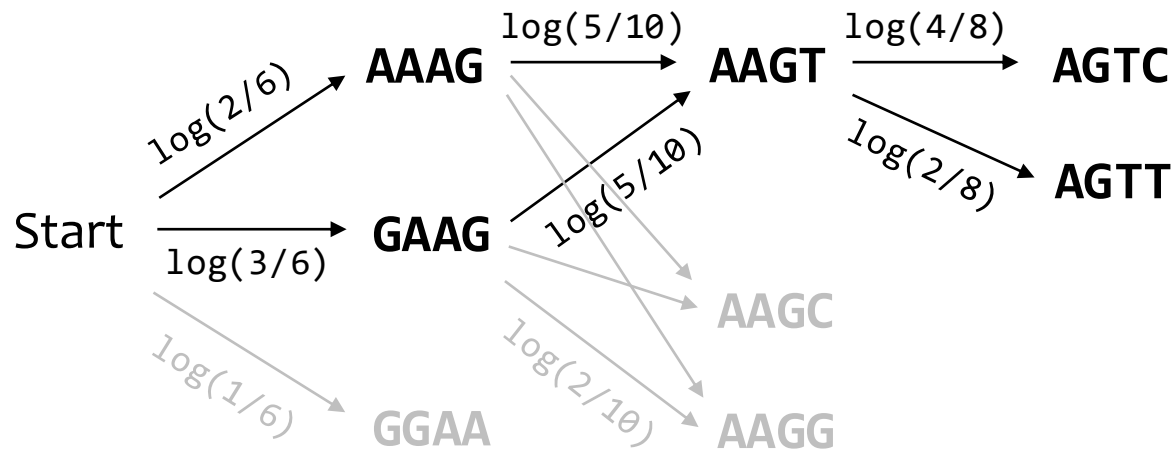
Goal: find the highest-weight path of length $L - k$ in the modified de Bruijn graph.



Beam search: keep only the top B paths (here $B = 2$).



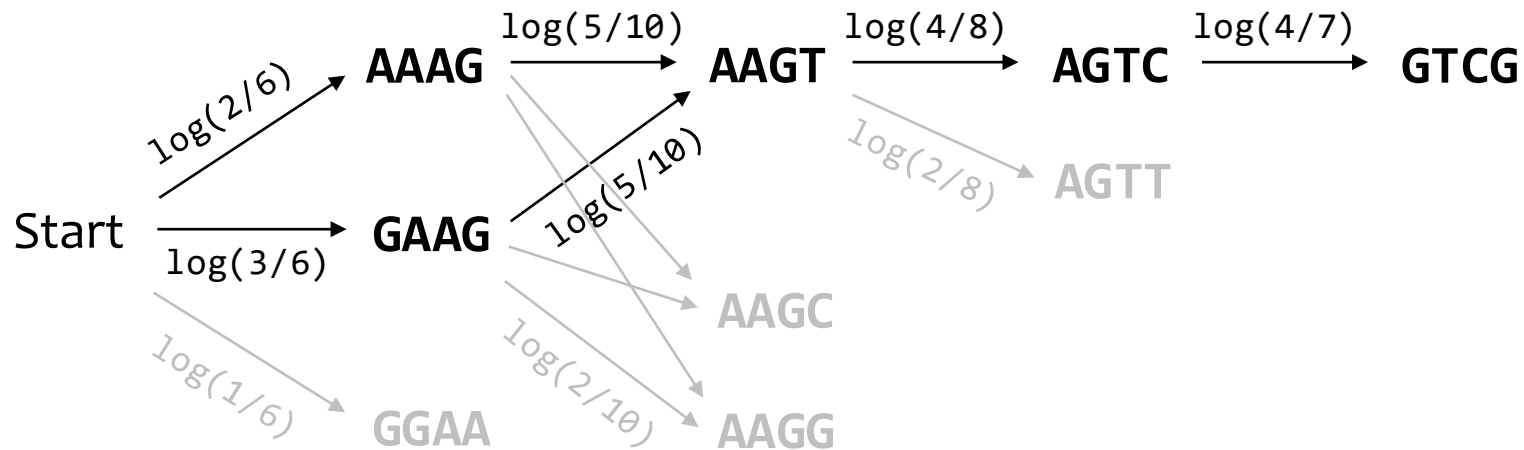
Goal: find the highest-weight path of length $L - k$ in the modified de Bruijn graph.



Beam search: keep only the top B paths (here $B = 2$).

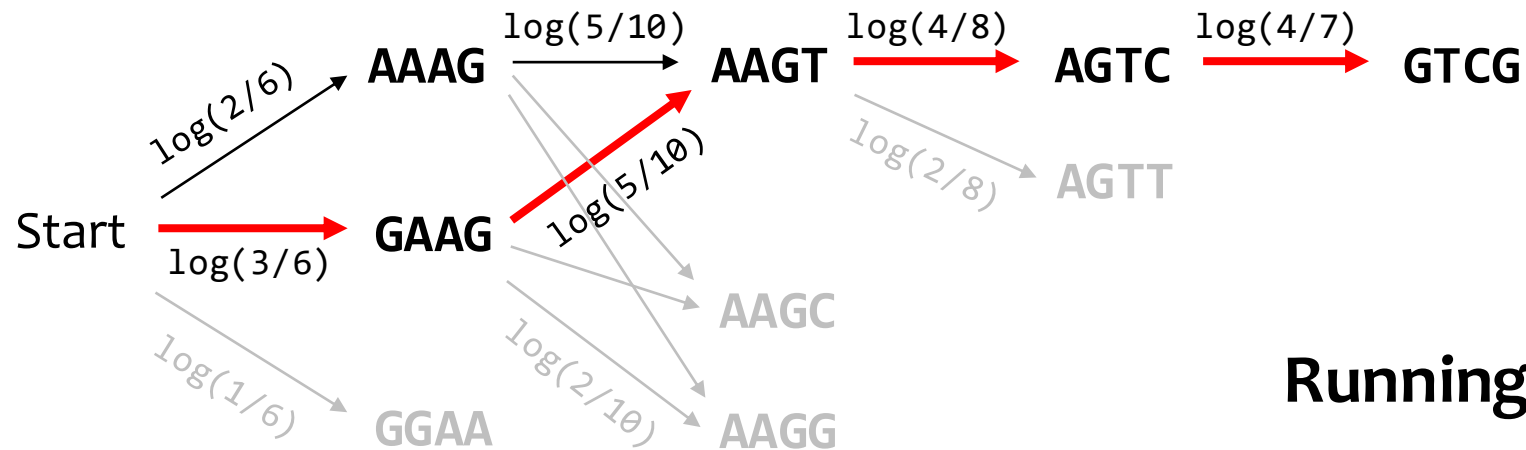


Goal: find the highest-weight path of length $L - k$ in the modified de Bruijn graph.



Beam search: keep only the top B paths (here $B = 2$).

Goal: find the highest-weight path of length $L - k$ in the modified de Bruijn graph.

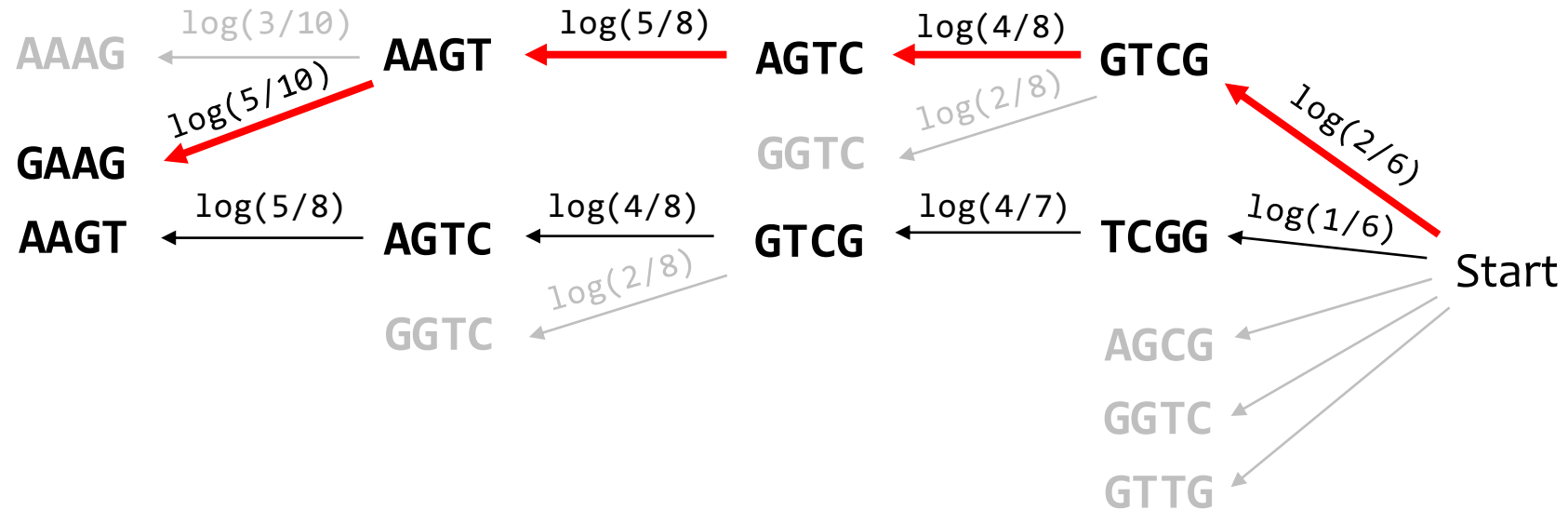


Running time: $O(BNL)$

Beam search: keep only the top B paths, $L - k$ iterations.



We can also do the same process from the end of the sequences.



Compare the two returned optimal paths by their weights.

Method

```
@002cd902-0726-4fb8-bf1d-0a056951d76f
GTGAGACTGACGACACAGCAGTGTAAGCAGTCCAAATGAAGCTTGGATGATTATTGACGTCTCCGTCGAATTCGTATCAGGGACACCCGTTTGTATGTGATGGACGAGAATGGATTCTCCTAGACCACAGACCCTGTGCGTGTGTGCGTGACTAGCTAGTACAT
+
+-4=>)JEIIBIOMQOS<RID@C)+( (&1?NGNNUUU+*&&' 253GMP)PKPK.=+9&,0320,,M8,(UUUUNI3>A*LRQLR)SOP)KU,-QP*K5,/.LO/JGCMTTPP@MQ,+08*%+0GPOKS)UTTRP4PPRTOL2E1L)PKL+QG)+KP-* .5+;5.'
@0039c1dc-ab1b-42fa-a601-59ac8ad4b60f
CTAGACGTGCGAGTATACTACTATGCCGAGTTATATTTAAACTGCATTTAGATTCCGTCGCTAGGCTACTGATGCCCCAGCTGCCGTTTGAATTAAGGAAGTGATACCTAGATGTGTGTTACAGCGACCTTAAATGACACGCAGTCAGACGCGTGT
+
**0,@)?EC=%+ALD4.-*HEJHB,,*)>M6POJRK+, ,JOGMGGCJ6, :+-K0B0F+.@-OKMME(' @+CF=)2'0/&FC+)QMKORT/PN3N+E7.. '564)RUQOJMMNIM.0/*)IRU*OMRNJ JNM+UQRQTA*88F-OULUT)MJO-56>5)%
@0041ac21-7266-4a7e-b9c8-1a1a92b729e2
TTAATGTAGCTATCTCTGCTAGCTAAATCGATCAGAAGTTTTCAAATCTTATAGTGATCTATAGCGCTTGATGGCGGAACGTGTCTTCGCTAAGTGCGGGCGTATGTGCAGCTTTGTGTCCCTGTATCGTTAACTAGGAAGCTGATGTACAGAGATAGACTACAG
+
&$&**6)(%+0,-/4+%(64;1+>OQ++JORK' '&)(3;-0&*,*#%'#$$/ '/1&&$*%--(*CB<MPENJKO+U50*%'&*49*(0')-+ ,/2HG6*1*;8+C93G)/.+)++)+(+%%, (&(%+( ('#,213D10Q0/+* :;/( *@OKNFB5'&$
```



```
@002cd902-0726-4fb8-bf1d-0a056951d76f READ IDENTIFIER
GTGAGACTGACGACACAGCAGTGTAAAGCAGTCCAAATGAAGCTTGGATGATTATTGACGTCTCCGTCGAATTCGTATCAGGGACACCCGTTTGGATGTGATGGACGAGAATGGATTCTCCTAGACCACAGACCCTGTGCGTGTGTGCGTACTAGCTAGTACAT
+
+-4=>)JEIIBIOMQOS<RID@C)+( (&1?NGNNUUU+*&&' 253GMP)PKPK.=+9&,0320,,M8,(UUUUNI3>A*LRQLR)SOP)KU,-QP*K5,/ .LO/JGCMTTPP@MQ,+08*%+0GPOKS)UTTRP4PPRTOL2E1L)PKL+QG)+KP-* .5+;5.'
@0039c1dc-ab1b-42fa-a601-59ac8ad4b60f
CTAGACGTGCGAGTATACTACTATGCCGAGTTATATTTAAACTGCATTTAGATTCCGTCGCTAGGCTACTGATGCCCCAGCTGCCGTTTGAATTAAGGAAGTGATACCTAGATGTGTGTTACAGCGACCTTAAATGACACGCAGTCAGACGCGTGT
+
**0,@)?EC=%+ALD4.-*HEJHB,,*)>M6POJRK+, ,JOGMGGCJ6, :+-K0B0F+.@-OKMME(' @+CF=)2'0/&FC+)QMKORT/PN3N+E7.. '564)RUQOJMMNIM.0/*)IRU*OMRNJNM+UQRQTA*88F-OULUT)MJO-56>5)%
@0041ac21-7266-4a7e-b9c8-1a1a92b729e2
TTAATGTAGCTATCTCTGCTAGCTAAATCGATCAGAAGTTTTCAAATCTTATAGTGATCTATAGCGCTTGATGGCGGAACGTGTCTTCGCTAAGTGCGGGCGTATGTGCAGCTTTGTGTCCCTGTATCGTTAACTAGGAAGCTGATGTACAGAGATAGACTACAG
+
&$&**6)(%%+0,-/4+%(64;1+>OQ++JORK' '&)(3;-%-0&*,*#%'#$$/ '/1&&$*%--(*CB<MPENJKO+U50*%'&*49*(0')-+ ,/2HG6*1*;8+C93G)/ .+)))+(+%%, (&(%+( ('#,213D10Q0/+* :;/( *@OKNFB5'&$
```



BASECALLED READ

```
@002cd902-0726-4fb8-bf1d-0a056951d76f
GTGAGACTGACGACACAGCAGTGTAAAGCAGTCCAAATGAAGCTTGGATGATTATTGACGTCTCCGTCGAATTCGTATCAGGGACACCCGTTTGTATGTGATGGACGAGAATGGATTCTCCTAGACCACAGACCCTGTGCGTGTGTGCGTACTAGCTAGTACAT
+
+-4=>)JEIIBIOMQOS<RID@C)+( (&1?NGNNUUU+*&&'253GMP)PKPK.=+9&,0320,,M8,(UUUUNI3>A*LRQLR)SOP)KU,-QP*K5,/.LO/JGCMTTPP@MQ,+08*%+0GPOKS)UTTRP4PPRTOL2E1L)PKL+QG)+KP-* .5+;5.'
@0039c1dc-ab1b-42fa-a601-59ac8ad4b60f
CTAGACGTGCGAGTATACTACTATGCCGAGTTATATTTAAACTGCATTTAGATTCCGTCGCTAGGCCTACTGATGCCCCAGCTGCCGTTTGAATTAAGGAAGTGATACCTAGATGTGTGTTACAGCGACCTTAAATGACACGCAGTCAGACGCGTGT
+
**0,@)?EC=%+ALD4.-*HEJHB,,*)>M6POJRK+, ,JOGMGGCJ6, :+-K0B0F+.@-OKMME(' @+CF=)2'0/&FC+)QMKORT/PN3N+E7.. '564)RUQOJMMNIM.0/*)IRU*OMRNJ JNM+UQRQTA*88F-OULUT)MJO-56>5)%
@0041ac21-7266-4a7e-b9c8-1a1a92b729e2
TTAATGTAGCTATCTCTGCTAGCTAAATCGATCAGAAGTTTTCAAATCTTATAGTGATCTATAGCGCTTGATGGCGGAACGTGTCTTCGCTAAGTGCGGGCGTATGTGCAGCTTTGTGTCCCTGTATCGTTAACTAGGAAGCTGATGTACAGAGATAGACTACAG
+
&$&**6)(%#+0,-/4+%(64;1+>OQ++JORK' '&)(3;-%-0&*,*#%'#$$/ '/1&&$*%--(*CB<MPENJKO+U50*%'&*49*(0')-+ ,/2HG6*1*;8+C93G)/.+)++)+(+%%, (&(%+( ('#,213D10Q0/+* :;/( *@OKNFB5'&$
```



```
@002cd902-0726-4fb8-bf1d-0a056951d76f
GTGAGACTGACGACACAGCAGTGTAAGCAGTCCAAATGAAGCTTGGATGATTATTGACGTCTCCGTCGAATTCTGATCAGGGACACCCGTTTGTATGTGATGGACGAGAATGGATTCTCCTAGACCACAGACCCTGTGCGTGTGTGCGTGACTAGCTAGTACAT
+ SEPARATOR
+-4=>)JEIIBIOMQOS<RID@C)+( (&1?NGNNUUU+*&&'253GMP)PKPK.=+9&,0320,,M8,(UUUUNI3>A*LRQLR)SOP)KU,-QP*K5,/ .LO/JGCMTPP@MQ,+08*%+0GPOKS)UTTRP4PPRTOL2E1L)PKL+QG)+KP-* .5+;5.'
@0039c1dc-ab1b-42fa-a601-59ac8ad4b60f
CTAGACGTGCGAGTATACTACTATGCCGAGTTATATTTAAACTGCATTTAGATTCCGTCGCTAGGCTACTGATGCCCCAGCTGCCGTTTGAATTAAGGAAGTGATACCTAGATGTGTGTTACAGCGACCTTAAATGACACGCAGTCAGACGCGTGT
+
**0,@)?EC=%+ALD4.-*HEJHB,,*)>M6POJRK+, ,JOGMGGCJ6, :+-K0B0F+.@-OKMME(' @+CF=)2'0/&FC+)QMKORT/PN3N+E7.. '564)RUQOJMMNIM.0/*)IRU*OMRNJ JNM+UQRQTA*88F-OULUT)MJO-56>5)%
@0041ac21-7266-4a7e-b9c8-1a1a92b729e2
TTAATGTAGCTATCTCTGCTAGCTAAATCGATCAGAAGTTTTCAAATCTTATAGTGATCTATAGCGCTTGATGGCGGAACGTGTCTTCGCTAAGTGCGGGCGTATGTGCAGCTTTGTGTCCCTGTATCGTTAACTAGGAAGCTGATGTACAGAGATAGACTACAG
+
&$&**6)(%%+0,-/4+%(64;1+>OQ++JORK' '&)(3;-%-0&*,*#%'#$$/ '/1&&$*%--(*CB<MPENJKO+U50*%'&*49*(0')-+,/2HG6*1*;8+C93G)/.+)++)+(+%%, (&(%+((' #,213D10Q0/+* :;/( *@OKNFB5'&$
```



QUALITY STRING

```
@002cd902-0726-4fb8-bf1d-0a056951d76f
GTGAGACTGACGACACAGCAGTGTAAAGCAGTCCAAATGAAGCTTGGATGATTATTGACGTCTCCGTCGAATTCTGATCAGGGACACCCGTTTGTATGTGATGGACGAGAATGGATTCTCCTAGACCACAGACCCTGTGCGTGTGTGCGTGACTAGCTAGTACAT
+
+-4=>)JEIIBIOMQOS<RID@C)+( (&1?NGNNUUU+*&&'253GMP)PKPK.=+9&,0320,,M8,(UUUUNI3>A*LRQLR)SOP)KU,-QP*K5,/ .LO/JGCMTTPP@MQ,+08*%+0GPOKS)UTTRP4PPRTOL2E1L)PKL+QG)+KP-* .5+;5.'
@0059c10c-a010-427a-a001-59ac8a04000f
CTAGACGTGCGAGTATACTACTATGCCGAGTTATATTTAAACTGCATTTAGATTCCGTCGCTAGGCTACTGATGCCCCAGCTGCCGTTTGAATTAAGGAAGTGATACCTAGATGTGTGTTTACAGCGACCTTAAATGACACGCAGTCAGACGCGTGT
+
**0,@)?EC=%+ALD4.-*HEJHB,,*) .>M6POJRK+, ,JOGMGGCJ6, :+-K0B0F+.@-OKMME(' @+CF=)2'0/&FC+)QMKORT/PN3N+E7.. '564)RUQOJMMNIM.0/*)IRU*OMRNJ JNM+UQRQTA*88F-OULUT)MJO-56>5)%
@0041ac21-7266-4a7e-b9c8-1a1a92b729e2
TTAATGTAGCTATCTCTGCTAGCTAAATCGATCAGAAGTTTTCAAATCTTATAGTGATCTATAGCGCTTGATGGCGGAACGTGTCTTCGCTAAGTGCGGGCGTATGTGCAGCTTTGTGTCCCTGTATCGTTAACTAGGAAGCTGATGTACAGAGATAGACTACAG
+
&$&**6)(%+0,-/4+(64;1+>OQ++JORK' '&)(3;-%-0&*,*#'#$$%/' /1&&$*%--(*CB<MPENJKO+U50*% '&*49*(0')-+,/2HG6*1*;8+C93G)/.+)++)+(+%%, (&(%+((' #,213D10Q0/+* :;//(*@OKNFB5'&$
```



Method

```
@002cd902-0726-4fb8-bf1d-0a056951d76f
GTGAGACTGACGACACAGCAGTGTAAAGCAGTCCAAATGAAGCTTGGATGATTATTGACGTCTCCGTC CAATTCTGATCAGGGACACCCGTTTGTATGTGATGGACGAGAATGGATTCTCCTAGACCACAGACCCTGTGCGTGTGTGCGTACTAGCTAGTACAT
+
+-4-;>)JEIIBIOMOQS<RID@C)+( (&1?NGNNUUU+*&&' 253GMP)PKPK.=+9&,0320,,M8,(UUUUNI3>A*LRQLR)SOP)KU,-QP*K5,/.LO/JGCMTTPP@MQ,+08*%+0GPOKS)UTTRP4PPRTOL2E1L)PKL+QG)+KP-*.5+;5.'
@0039c1dc-ab1b-42fa-a601-59ac8ad4b60f
CTAGACGTGCGAGTATACTACTATGCCGAGTTATATTTAAACTGCATTTAGATTCCGTCGCTAGGCTACTGATGCCCCAGCTGCCGTTTGAATTAAGGAAGTGATACCTAGATGTGTGTTACAGCGACCTTAAATGACACGCAGTCAGACGCGTGT
+
**0,@)?EC=%+ALD4.-*HEJHB,,*)>M6POJRK+, ,JOGMGGCJ6, :+-K0B0F+.@-OKMME(' @+CF=)2'0/&FC+)QMKORT/PN3N+E7.. '564)RUQOJMMNIM.0/*)IRU*OMRNJ JNM+UQRQTA*88F-OULUT)MJO-56>5)%
@0041ac21-7266-4a7e-b9c8-1a1a92b729e2
TTAATGTAGCTATCTCTGCTAGCTAAATCGATCAGAAGTTTTCAAATCTTATAGTGATCTATAGCGCTTGATGGCGGAACGTGTCTTCGCTAAGTGCGGGCGTATGTGCAGCTTTGTGTCCCTGTATCGTTAACTAGGAAGCTGATGTACAGAGATAGACTACAG
+
&$&**6)(%#+0,-/4+%(64;1+>OQ++JORK' '&)(3;-%-0&*,*#%'#$$/ '/1&&$*%--(*CB<MPENJKO+U50*%'&*49*(0')-+/,2HG6*1*;8+C93G)/.+)++)+(+%%, (&(%+((' #,213D10Q0/+* :;/( *@OKNFB5'&$&
```

Basecalled Base	ASCII Quality Score	Q-Score	P(error)	P(correct)
G	+			
T	-			
G	4			



```
@002cd902-0726-4fb8-bf1d-0a056951d76f
GTG GACTGACGACACAGCAGTGTAAAGCAGTCCAAATGAAGCTTGGATGATTATTGACGTCTCCGTCGAATTCGTATCAGGGACACCCGTTTGTATGTGATGGACGAGAATGGATTCTCCTAGACCACAGACCCTGTGCGTGTGTGCGTGTACTAGCTAGTACAT
+
+-4-:>)JEIIBIOMQOS<RID@C)+( (&1?NGNNUUU+*&&' 253GMP)PKPK.=+9&,0320,,M8,(UUUUNI3>A*LRQLR)SOP)KU,-QP*K5,/.LO/JGCMTPP@MQ,+08*%+0GPOKS)UTTRP4PPRTOL2E1L)PKL+QG)+KP-*.5+;5.'
@0039c1dc-ab1b-42fa-a601-59ac8ad4b60f
CTAGACGTGCGAGTATACTACTATGCGGAGTTATATTTAAACTGCATTTAGATTCCGTCGCTAGGCTACTGATGCCCCAGCTGCCGTTTGAATTAAGGAAGTGATACCTAGATGTGTGTTACAGCGACCTTAAATGACACGCAGTCAGACGCGTGT
+
**0,@)?EC=%+ALD4.-*HEJHB,,*)>M6POJRK+, ,JOGMGGCJ6, :+-K0B0F+.@-OKMME('@+CF=)2'0/&FC+)QMKORT/PN3N+E7.. '564)RUQOJMMNIM.0/*)IRU*OMRNJ JNM+UQRQTA*88F-OULUT)MJO-56>5)%
@0041ac21-7266-4a7e-b9c8-1a1a92b729e2
TTAATGTAGCTATCTCTGCTAGCTAAATCGATCAGAAGTTTTCAAATCTTATAGTGATCTATAGCGCTTGATGGCGGAACGTGTCTTCGCTAAGTGCGGGCGTATGTGCAGCTTTGTGTCCCTGTATCGTTAACTAGGAAGCTGATGTACAGAGATAGACTACAG
+
&$&**6)(%#+0,-/4+%(64;1+>OQ++JORK' '&)(3;-%-0&*,*#%'#$$/ '/1&&$*%--(*CB<MPENJKO+U50*%'&*49*(0')-+/,2HG6*1*;8+C93G)/.+)++)+(+%%, (&(%+((' #,213D10Q0/+*://(*@OKNFB5'&$&
```

$$Q\text{-score} = \text{ASCII value} - 33$$

$$P(\text{error}) = 10^{-Q/10}$$

$$P(\text{correct}) = 1 - P(\text{error})$$

Basecalled Base	ASCII Quality Score	Q-Score	P(error)	P(correct)
G	+			
T	-			
G	4			



```
@002cd902-0726-4fb8-bf1d-0a056951d76f
GTGAGACTGACGACACAGCAGTGTAAAGCAGTCCAAATGAAGCTTGGATGATTATTGACGTCTCCGTCCTCAATTCGTATCAGGGACACCCGTTTGGATGTGATGGACGAGAATGGATTCTCCTAGACCACAGACCCTGTGCGTGTGTGCGTGTACTAGCTAGTACAT
+
+-4-;>)JEIIBIOMQOS<RID@C)+( (&1?NGNNUUU+*&&'253GMP)PKPK.=+9&,0320,,M8,(UUUUNI3>A*LRQLR)SOP)KU,-QP*K5,/.LO/JGCMTPP@MQ,+08*%+0GPOKS)UTTRP4PPRTOL2E1L)PKL+QG)+KP-*.5+;5.'
@0039c1dc-ab1b-42fa-a601-59ac8ad4b60f
CTAGACGTGCGAGTATACTACTATGCGGAGTTATATTTAAACTGCATTTAGATTCCGTCGCTAGGCTACTGATGCCCCAGCTGCCGTTTGAATTAAGGAAGTGATACCTAGATGTGTGTTACAGCGACCTTAAATGACACGCAGTCAGACGCGTGT
+
**0,@)?EC=%+ALD4.-*HEJHB,,*)>M6POJRK+, ,JOGMGGCJ6, :+-K0B0F+.@-OKMME('@+CF=)2'0/&FC+)QMKORT/PN3N+E7.. '564)RUQOJMMNIM.0/*)IRU*OMRNJNM+UQRQTA*88F-OULUT)MJO-56>5)%
@0041ac21-7266-4a7e-b9c8-1a1a92b729e2
TTAATGTAGCTATCTCTGCTAGCTAAATCGATCAGAAGTTTTCAAATCTTATAGTGATCTATAGCGCTTGATGGCGGAACGTGTCTTCGCTAAGTGCGGGCGTATGTGCAGCTTTGTGTCCCTGTATCGTTAACTAGGAAGCTGATGTACAGAGATAGACTACAG
+
&$&**6)(%+0,-/4+(64;1+>OQ++JORK' '&)(3;-%-0&*,*#%'#$$/ '/1&&$*%--(*CB<MPENJKO+U50*%'&*49*(0')-+/,2HG6*1*;8+C93G)/.+)++)+(+%%, (&(%+((' #,213D10Q0/+*://(*@OKNFB5'&$
```

$$Q\text{-score} = \text{ASCII value} - 33$$

$$P(\text{error}) = 10^{-Q/10}$$

$$P(\text{correct}) = 1 - P(\text{error})$$

Basecalled Base	ASCII Quality Score	Q-Score	P(error)	P(correct)
G	+	10	0.1000	0.9000
T	-	12	0.0631	0.9369
G	4	19	0.0126	0.9874



BBS: count-only update

Observed k-mer: AGTC

```
count(AGTC) += 1
```



BBS: count-only update

Observed k-mer: AGTC

count(AGTC) += 1



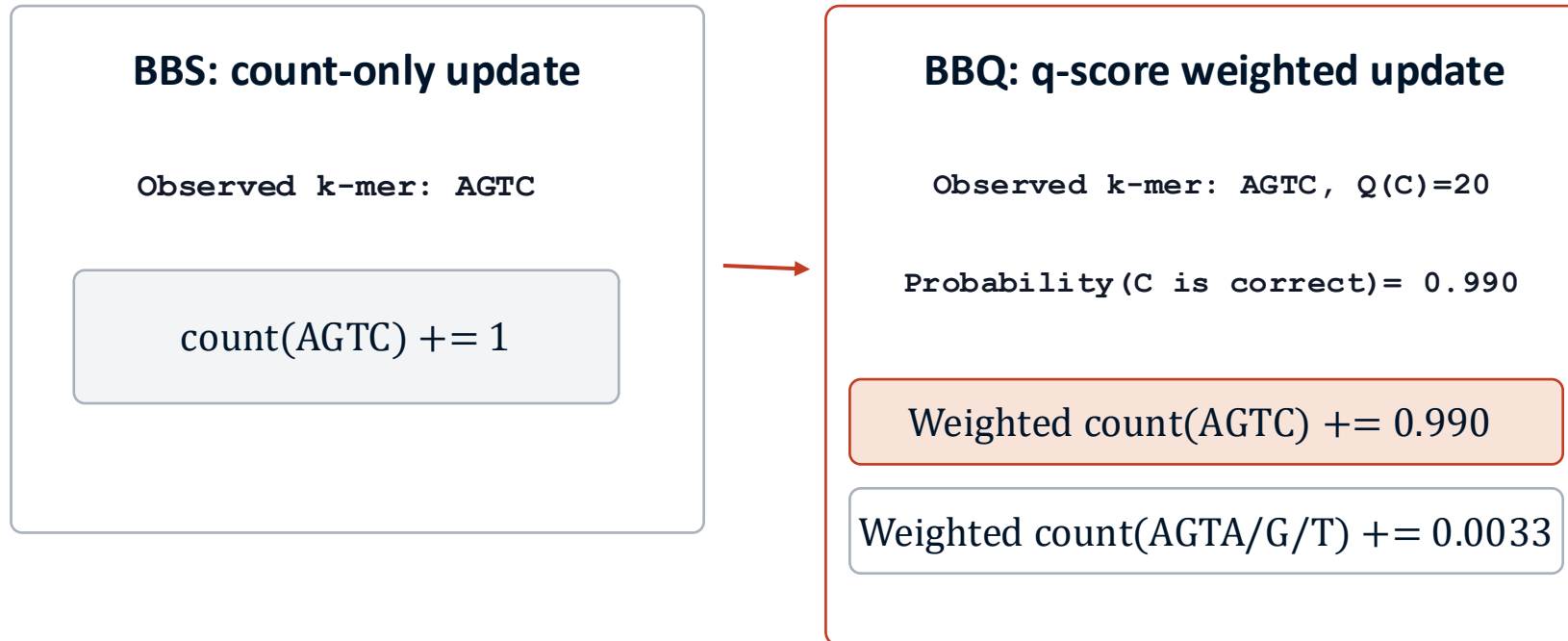
BBQ: q-score weighted update

Observed k-mer: AGTC, $Q(C)=20$

Probability(C is correct) = 0.990

Weighted count(AGTC) += 0.990

Weighted count(AGTA/G/T) += 0.0033



- k is chosen exactly as in BBS.
- The same pseudo-count smoothing can be used.
- Each transition still means: keep the prefix and append one base.

Results



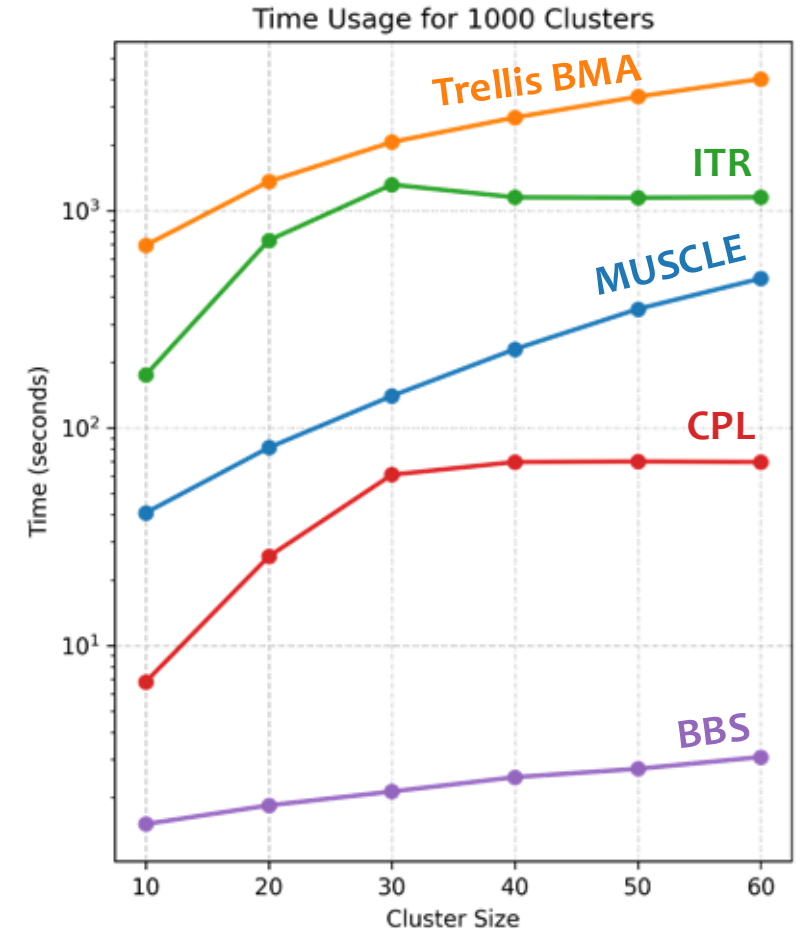
Real datasets from different sources.

	Bar-Lev et al., 2025	Srinivasavaradhan et al., 2021	Chandak et al., 2020
# Clusters	10000	9984	1466
Synthesis Technology	Twist Bioscience	Twist Bioscience	CustomArray
Sequencing Technology	MinION	MinION	MinION
Clustering algorithm	[Bar-Lev et al., 2025]	[Rashtchian et al., 2018]	Perfect
L	140	110	108
Coverage	21.37	27.01	114.29
Read error rate	4.34%	5.77%	13.56%



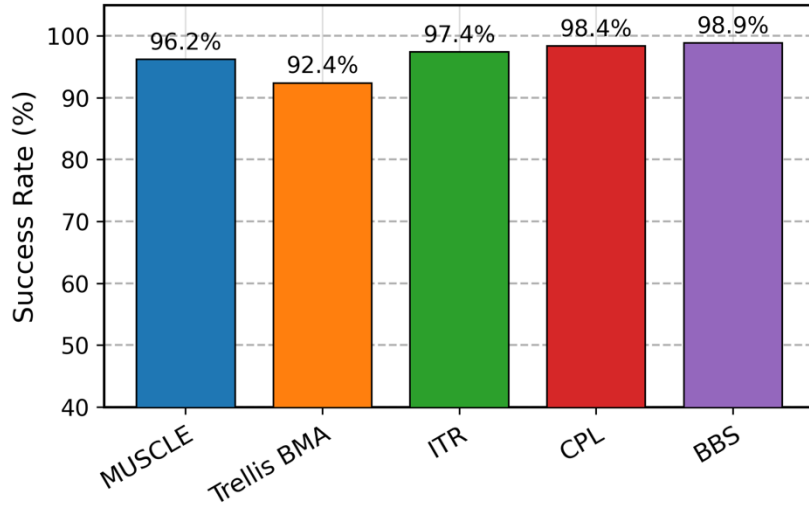
Results

Algorithm	Method
MUSCLE [Edgar, 2004] [Antkowiak et al., 2020]	MSA
Trellis BMA [Srinivasavaradhan et al., 2021]	Read correction
ITR [Sabary et al., 2024]	Assembly
CPL [Bar-Lev et al., 2025]	Pairwise alignment
BBS & BBQ (this work, beam width $B = 20$)	k-th order Markov chain

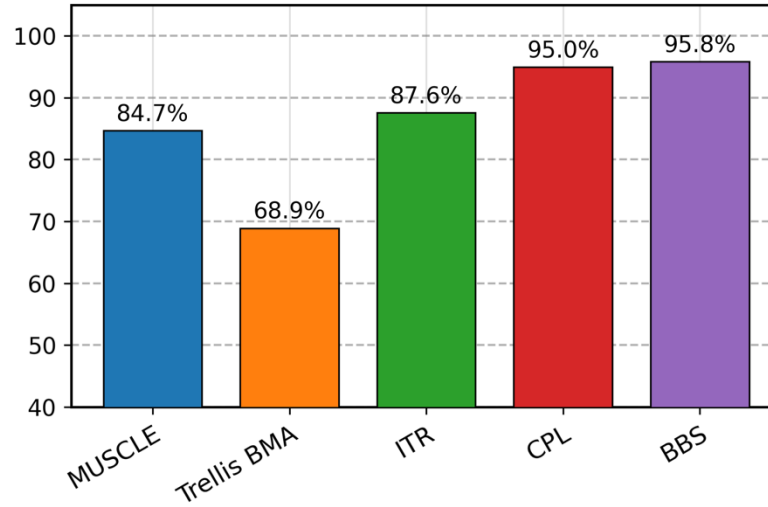


Results

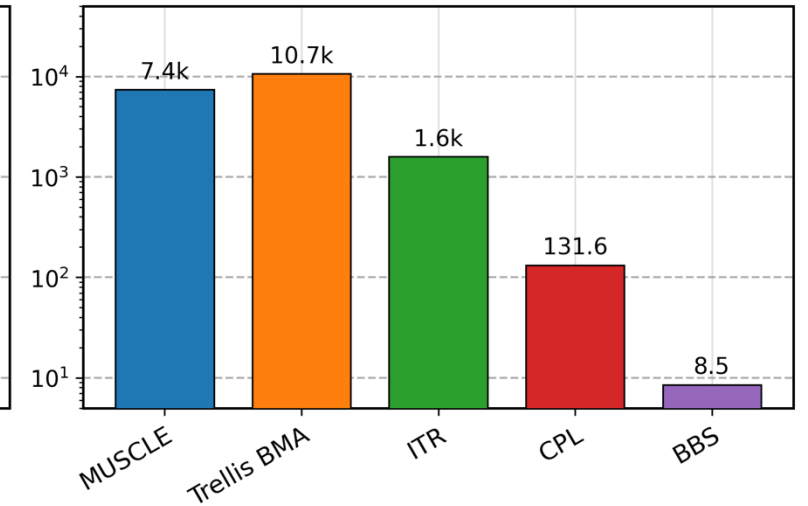
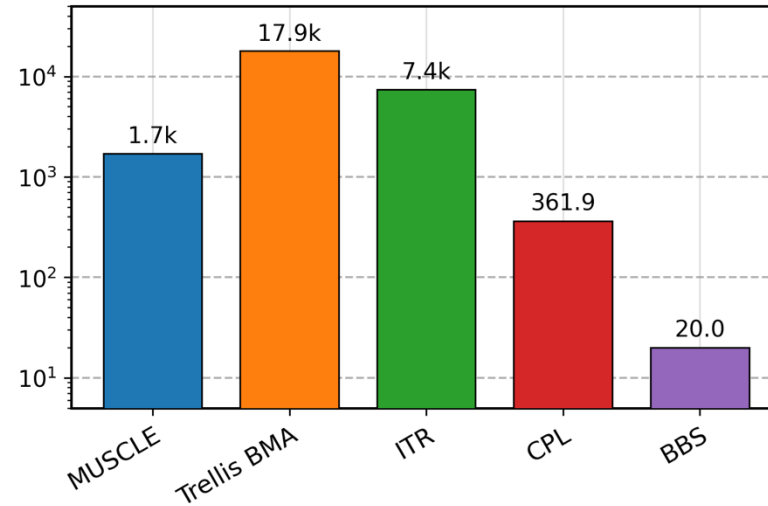
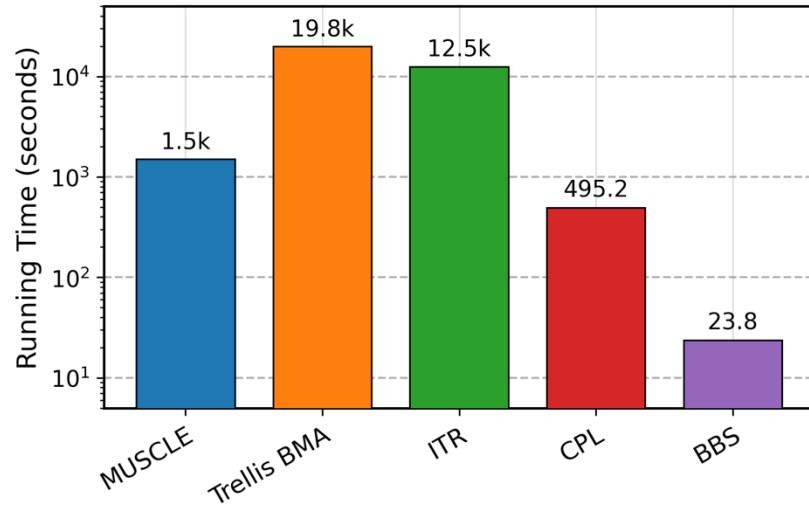
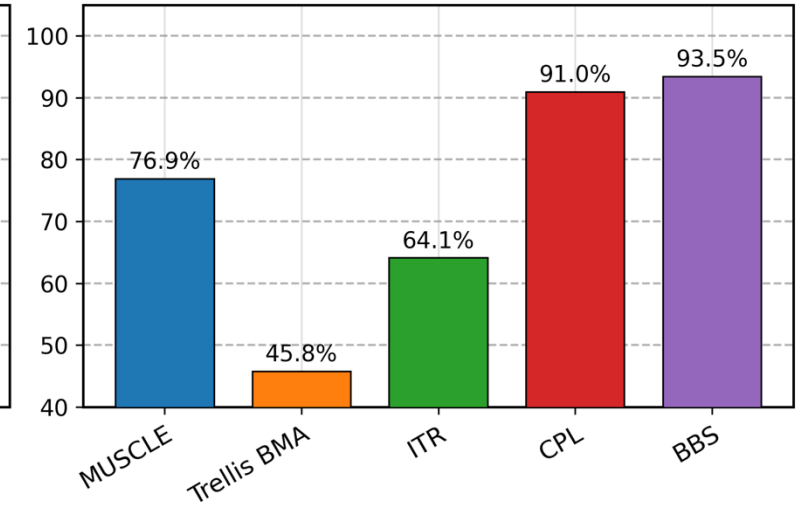
Bar-Lev et al.
(10k clusters, error rate 4.34%)



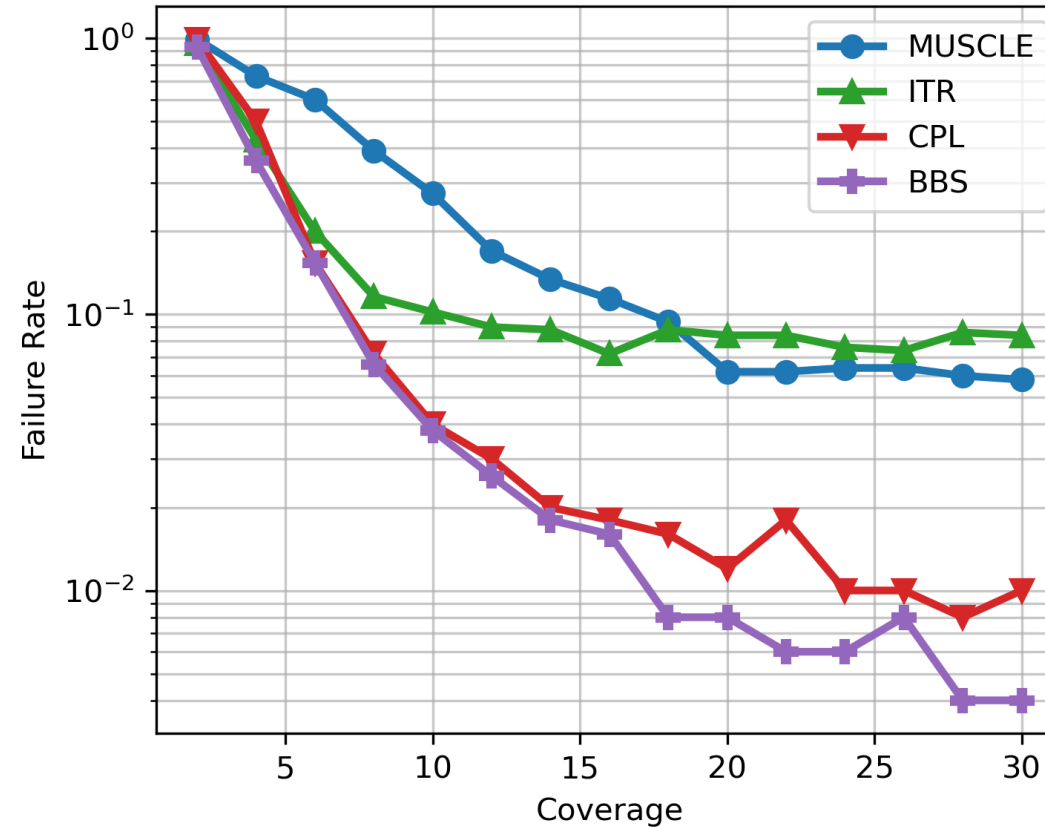
Srinivasavaradhan et al.
(10k clusters, error rate 5.77%)



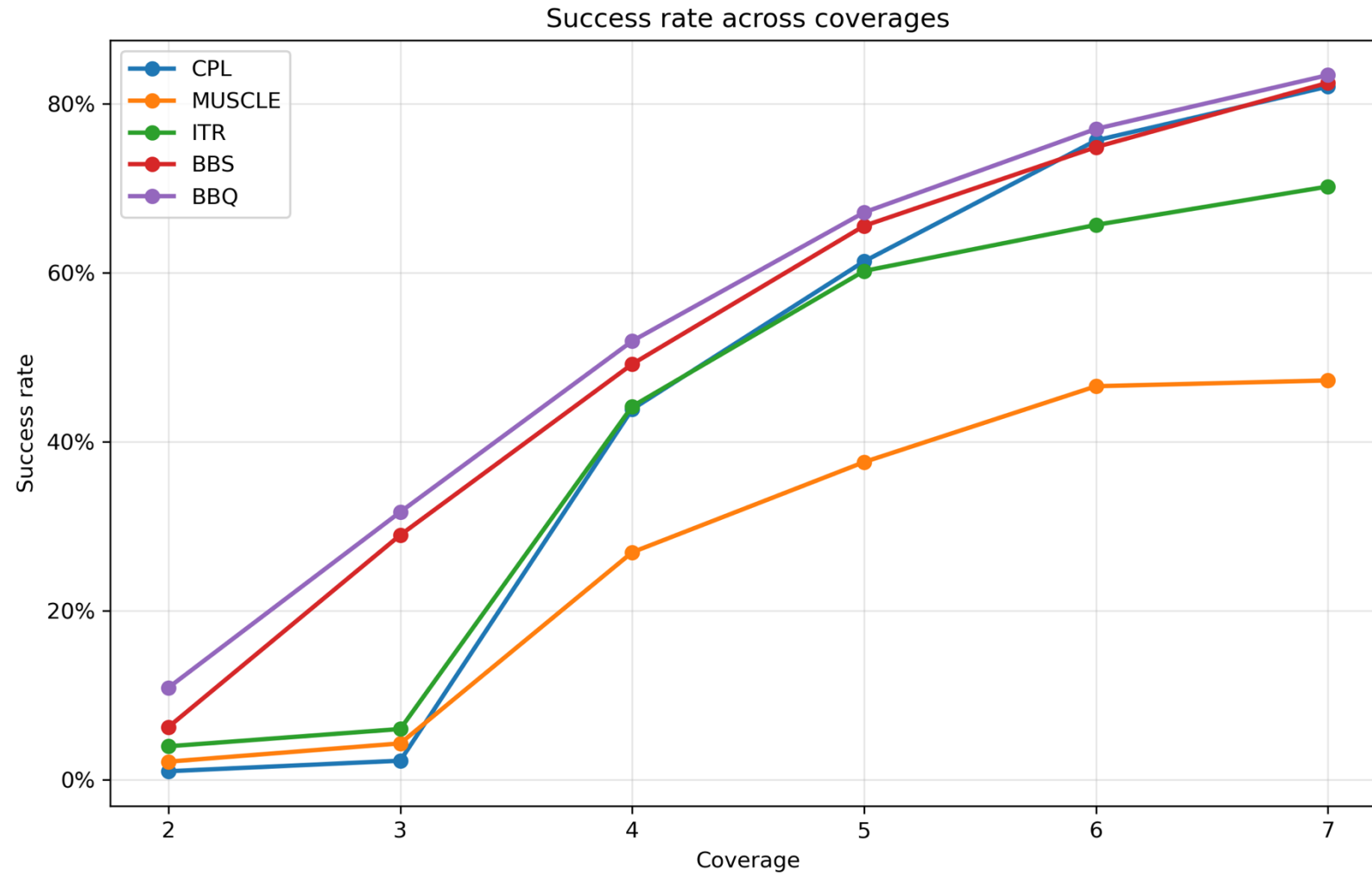
Chandak et al.
(1.5k clusters, error rate 13.56%)

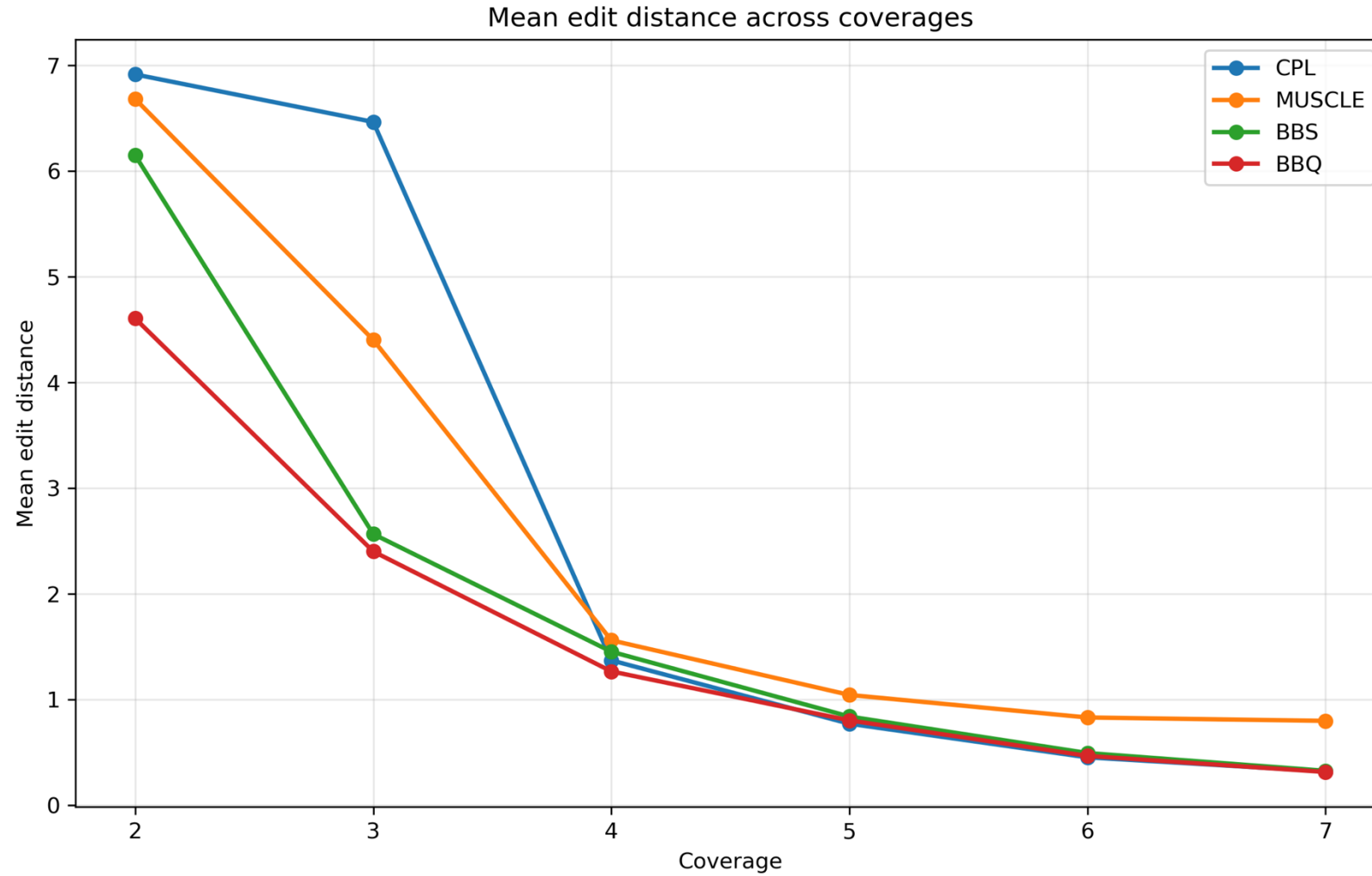


(Lower is better)



BBS performs well even with small coverage.





Takeaways:

- Modeling the reads as outputs of a k -th order Markov chain is an effective way to capture position-specific error rates.
- Beam search is very fast compared to traditional methods based on alignment or assembly, and has almost no loss in accuracy.
- The path weights of returned candidates are good indicators of reconstruction quality.
- The q-score in FASTQ files is an important tool to achieve higher reconstruction accuracy at low coverages.



Acknowledgement

Thank you very much for your attention!

My collaborators:



Zhenhao
Gu



Gary
Goh Yipeng



Weng-Fai
Wong

Special thanks:

Djordje Jevdic , Vasily Shenshin, Yannick Rondelez

STORAGE AND COMPUTING WITH **Dna**

Thank you

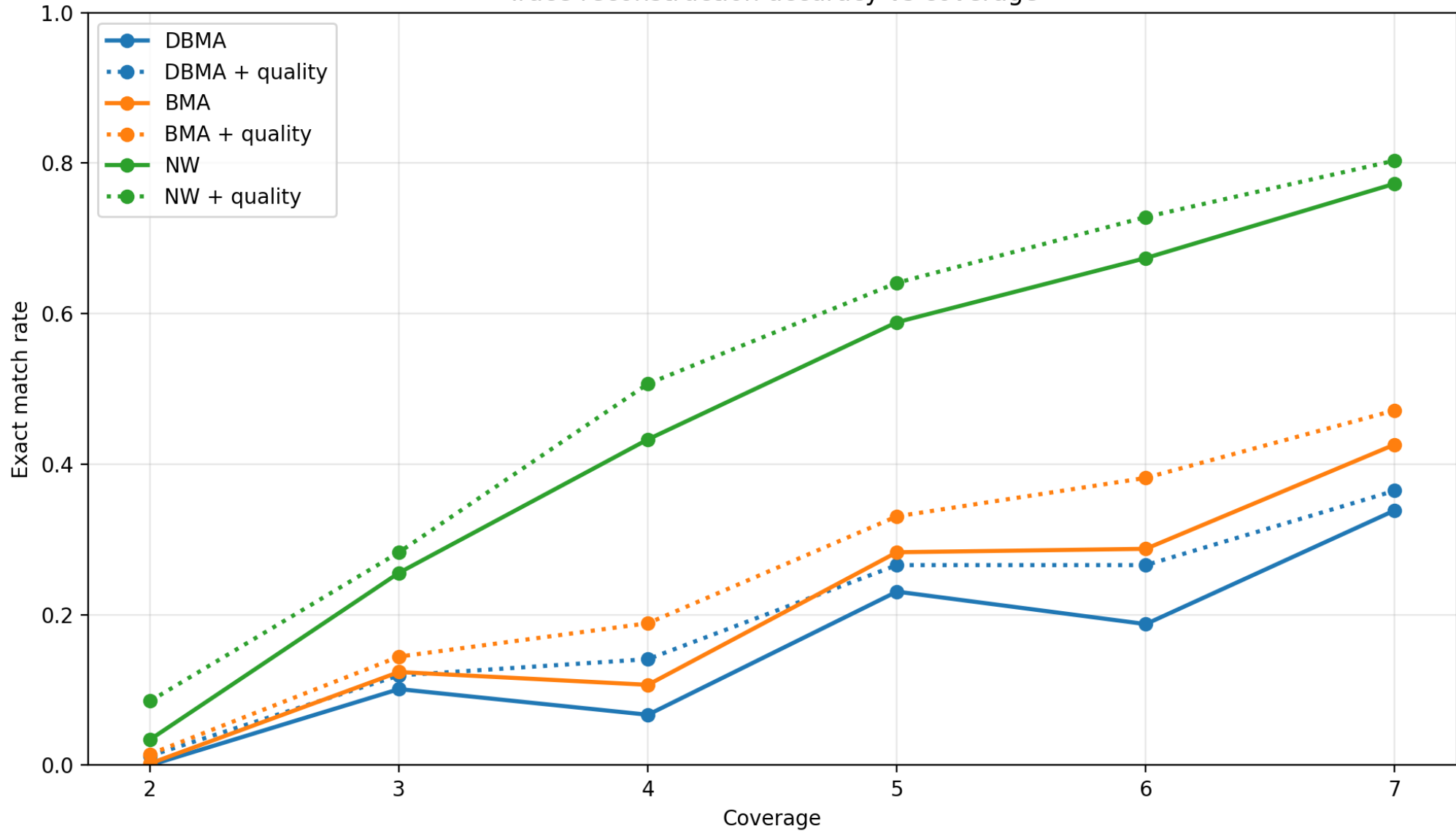
Sponsored by



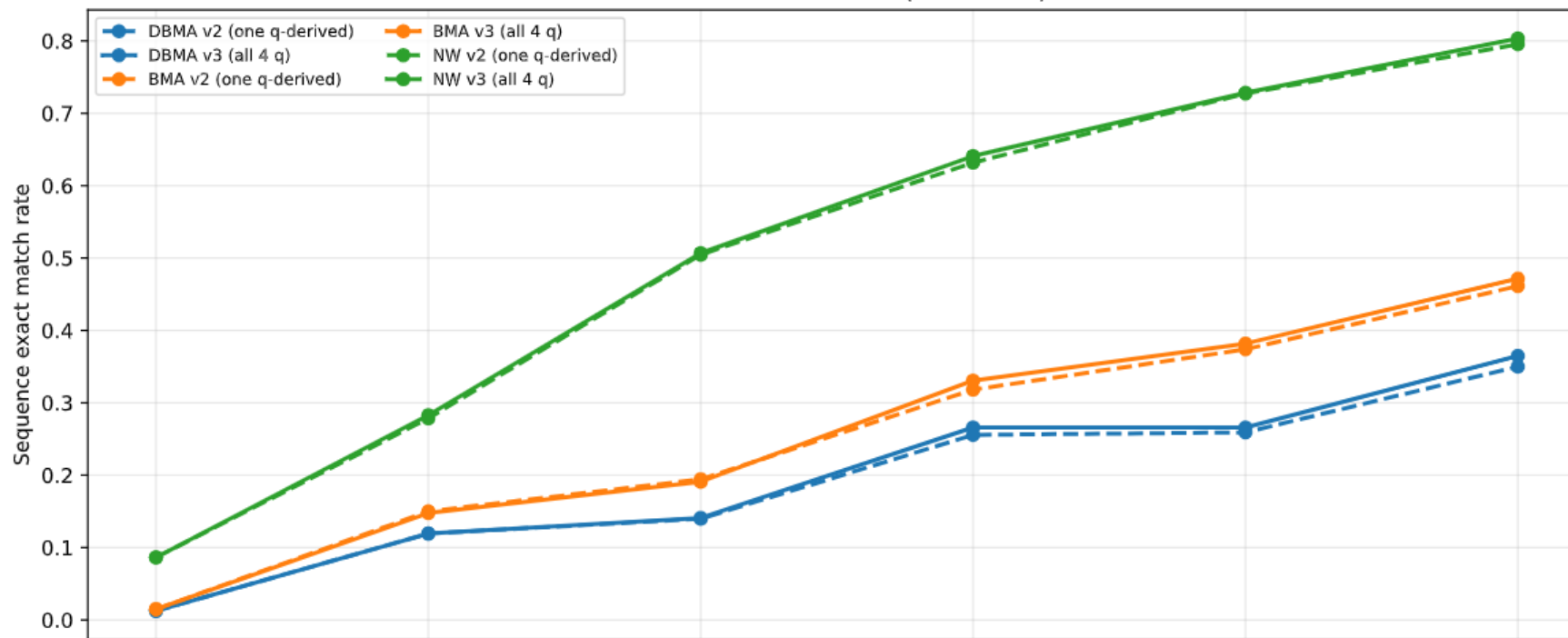
Organized by



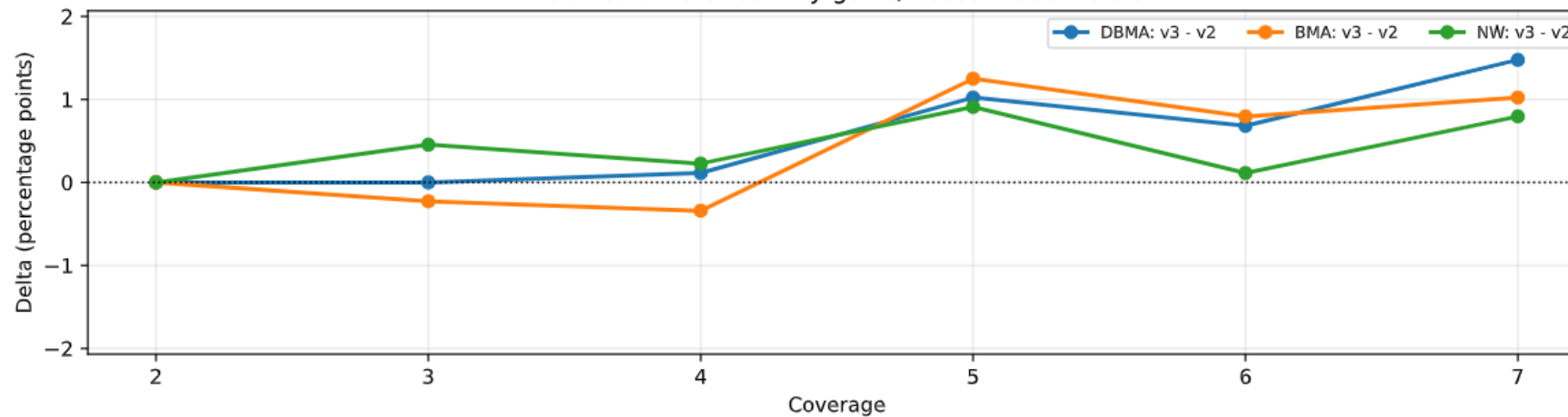
Trace reconstruction accuracy vs coverage



Version 2 vs Version 3 (main view)



Zoomed difference: tiny gains/losses made visible



b3s: Called-only vs Full A/C/G/T q-score weighting (zoomed)

