

Bit-encoding of canonical k -mers

Roland Wittler

Bielefeld University

DSB 2022

k -mers = sequences of length k over alphabet $\{A, C, G, T\}$

ATCGACTAGCTACAGCGT...

...

k -mers

k -mers = sequences of length k over alphabet $\{A, C, G, T\}$

ATCGACTAGCTACAGCGT...

...

Encoding

$$4^k \begin{cases} \text{AAAAA} \\ \text{AAAAC} \\ \vdots \\ \text{TTTTT} \end{cases} \Rightarrow \begin{matrix} \text{A} \mapsto 00 \\ \text{C} \mapsto 01 \\ \text{G} \mapsto 10 \\ \text{T} \mapsto 11 \end{matrix} \Rightarrow \begin{matrix} 00 & 00 & 00 & 00 & 00 \\ 00 & 00 & 00 & 00 & 01 \\ & & & \vdots & \\ 11 & 11 & 11 & 11 & 11 \end{matrix} \Rightarrow \begin{matrix} 0 \\ 1 \\ \vdots \\ 4^k - 1 \end{matrix}$$

Canonical k -mers (k odd)

Canonical representative of k -mer $x = \min_{lex} \{x, \bar{x}\}$

AAAAA	AAAAC	...	ATTTT	...	TAAAA	...	TTTTG	TTTTT
TTTTT	GTTTT		AAAAT		TTTTA		CAAAA	AAAAA

Canonical k -mers (k odd)

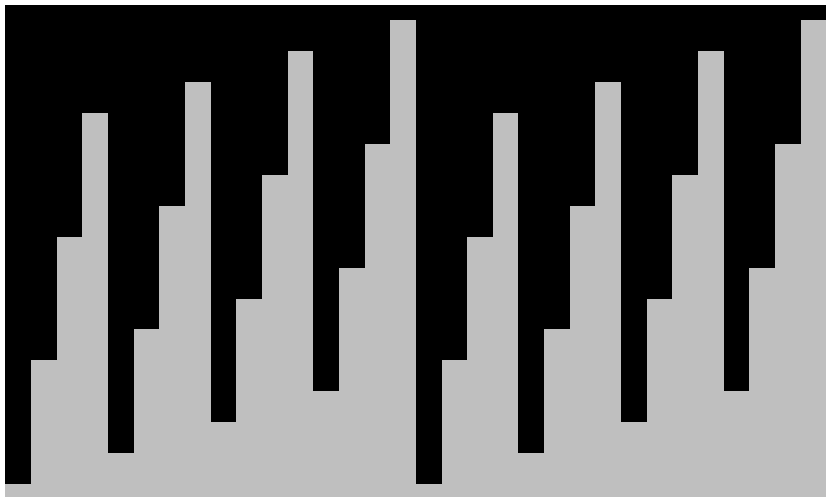
Canonical representative of k -mer $x = \min_{lex} \{x, \bar{x}\}$

AAAAA AAAAC ... ATTTT ... TAAAA ... TTTTG TTTTT
TTTTT GTTTT ... AAAAT ... TTTTA ... CAAAA AAAAA

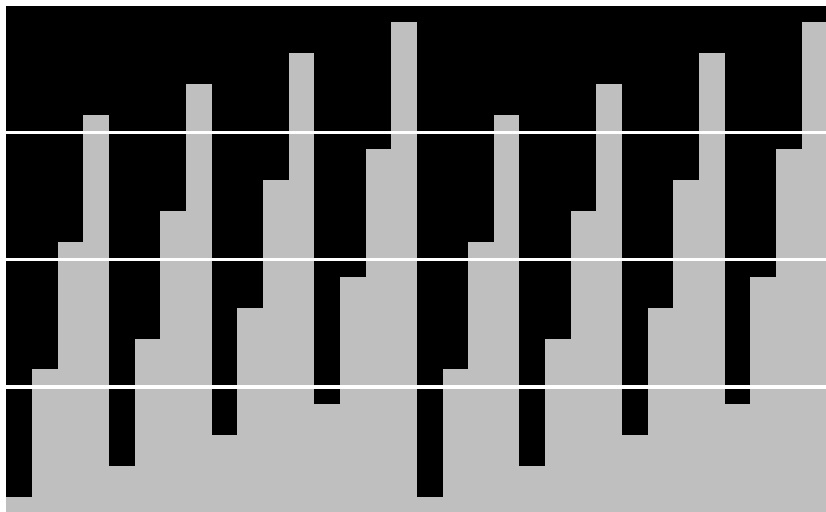
Encoding

$\frac{4^k}{2} \left\{ \begin{array}{l} \text{AAAAA} \\ \text{AAAAC} \\ \vdots \\ \text{TTTTT} \end{array} \right. \Rightarrow \begin{array}{l} \text{A} \mapsto 00 \\ \text{C} \mapsto 01 \\ \text{G} \mapsto 10 \\ \text{T} \mapsto 11 \end{array} \Rightarrow \begin{array}{l} 00 \ 00 \ 00 \ 00 \ 00 \\ 00 \ 00 \ 00 \ 00 \ 01 \\ \vdots \\ 11 \ 11 \ 11 \ 11 \ 11 \end{array} \Rightarrow \begin{array}{l} 0 \\ 1 \\ \vdots \\ 4^k - 1 \end{array}$

Canonical k -mers (k odd)



Canonical k -mers (k odd)

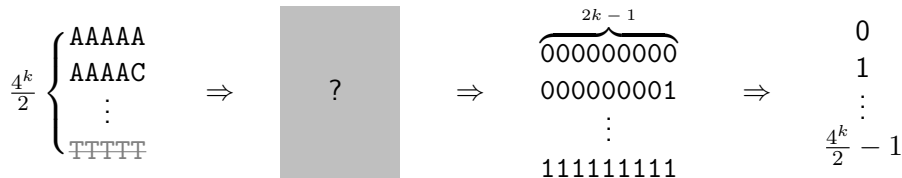


New encoding

canonical k -mers = $\frac{1}{2}$ # k -mers

► $2k - 1$ bits per k -mer

Encoding



How to save one bit?

- ▶ What makes a k -mer canonical?

New encoding

How to save one bit?

- ▶ What makes a k -mer canonical?

Encode from the outside inwards

A G A T G A T

New encoding

How to save one bit?

- ▶ What makes a k -mer canonical?

Encode from the outside inwards

A G A T G A **T**

New encoding

How to save one bit?

- ▶ What makes a k -mer canonical?

Encode from the outside inwards

A G A T G A T

New encoding (k odd)

Encode from the outside inwards

A G A T G A T
11 11

unspecific pair

A ... T \mapsto 11 ... 11

C ... G \mapsto 11 ... 10

G ... C \mapsto 11 ... 01

T ... A \mapsto 11 ... 00

New encoding (k odd)

Encode from the outside inwards

A G A T G A T
11 10 1 11

unspecific pair

A...T \mapsto 11...11

C...G \mapsto 11...10

G...C \mapsto 11...01

T...A \mapsto 11...00

specifying case

A...A \mapsto 00...0

A...C \mapsto 00...1

A...G \mapsto 01...0

C...A \mapsto 01...1

C...C \mapsto 10...0

G...A \mapsto 10...1

A \mapsto 0

C \mapsto 1

New encoding (k odd)

Encode from the outside inwards

A G A T G A T
11 10 00 11 10 1 11

unspecific pair

A...T \mapsto 11...11

C...G \mapsto 11...10

G...C \mapsto 11...01

T...A \mapsto 11...00

specifying case

A...A \mapsto 00...0

A...C \mapsto 00...1

A...G \mapsto 01...0

C...A \mapsto 01...1

C...C \mapsto 10...0

G...A \mapsto 10...1

A \mapsto 0

C \mapsto 1

remainder

A \mapsto 00

C \mapsto 01

G \mapsto 10

T \mapsto 11

Even-length k -mers

Palindromes

- ▶ $x = \bar{x}$ = canonical
- ▶ # palindromes = $4^{k/2}$
- ▶ # non-palindromes = $4^k - 4^{k/2}$
- ▶ # canonical k -mers = $\frac{1}{2}4^k - \frac{1}{2}4^{k/2} + 4^{k/2} > \frac{1}{2}4^k$
- ▶ more than $2k - 1$ bits $\Rightarrow 2k$ bits

Even-length k -mers

Palindromes

- ▶ $x = \bar{x}$ = canonical
- ▶ # palindromes = $4^{k/2}$
- ▶ # non-palindromes = $4^k - 4^{k/2}$
- ▶ # canonical k -mers = $\frac{1}{2}4^k - \frac{1}{2}4^{k/2} + 4^{k/2} > \frac{1}{2}4^k$
- ▶ more than $2k - 1$ bits $\Rightarrow 2k$ bits

Tweak new encoding

- ▶ span integer range $0, \dots, \# \text{ canonical } k\text{-mers} - 1$
- ▶ palindromes get highest ranks

Even-length k -mers

Palindromes get highest ranks?!

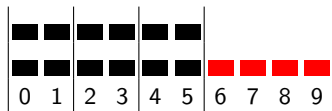
- ▶ non-palindromes: two k -mers = one canonical k -mer
- ▶ palindrome: one k -mer = one canonical k -mer



Even-length k -mers

Palindromes get highest ranks?!

- ▶ non-palindromes: two k -mers = one canonical k -mer
- ▶ palindrome: one k -mer = one canonical k -mer
- ▶ distribute canonical k -mers to buckets: heterogeneous bucket sizes



Even-length k -mers

0000000	AAAA	01101110	CACG	01111000	TCGA
0000001	AAAC	01101110	CCAG	01111001	GCGC
0000010	AACA	01101111	AACT	01111010	CCGG
0000011	AACC	01101111	ACAT	01111011	ACGT
		01110000	TCCA	01111100	TATA
		01110001	GCCC	01111101	GATC
		01110010	CCCG	01111110	CATG
		01110011	ACCT	01111111	AATT
		01110100	TAGA	10000000	TTAA
		01110100	TGAA	10000001	GTAC
		01110101	GAGC	10000010	CTAG
		01110101	GGAC	10000011	ATAT
01101100	TACA	01110110	CAGG	10000100	TGCA
01101100	TCAA	01110110	CGAG	10000101	GGCC
01101101	GACC	01110111	AAGT	10000110	CGCG
01101101	GCAC	01110111	AGAT	10000111	AGCT

Proof of concept

Frigate [Brihadiswaran & Jayasena, Int. Conf. Bioinf. & Biomed. Techn., 2021]

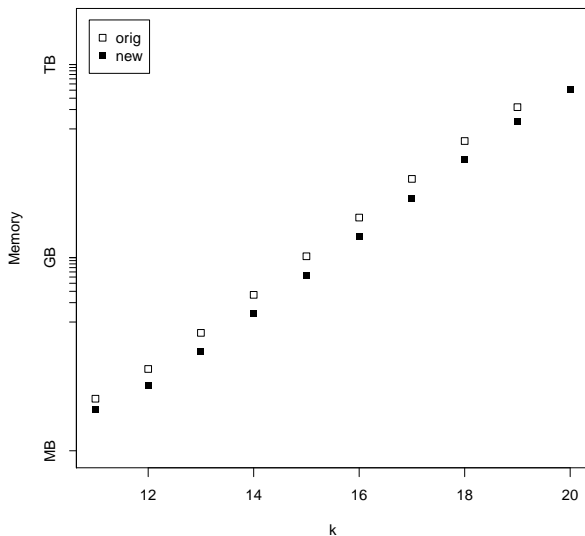
- ▶ k -mer counter
- ▶ in-memory: array of length 4^k
- ▶ multithreaded, lock-free array access

```
minimum = kmer;  
if ( rc_kmer < kmer ) {  
    minimum = rc_kmer;  
}
```

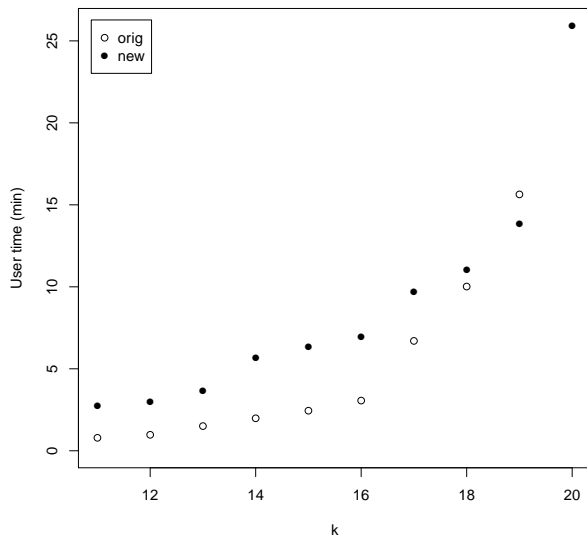
⇒ `minimum = new_enc(kmer, rc_kmer); //const. time`

```
pow(4, k_value)  
⇒ pow(4, k_value) / 2 //k odd
```

F. vesca

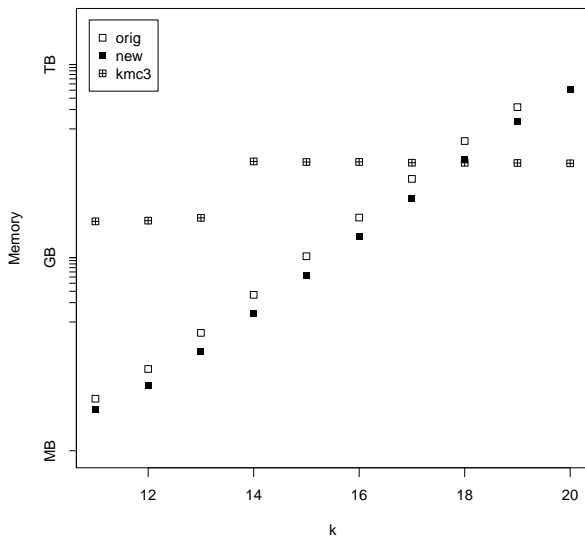


F. vesca

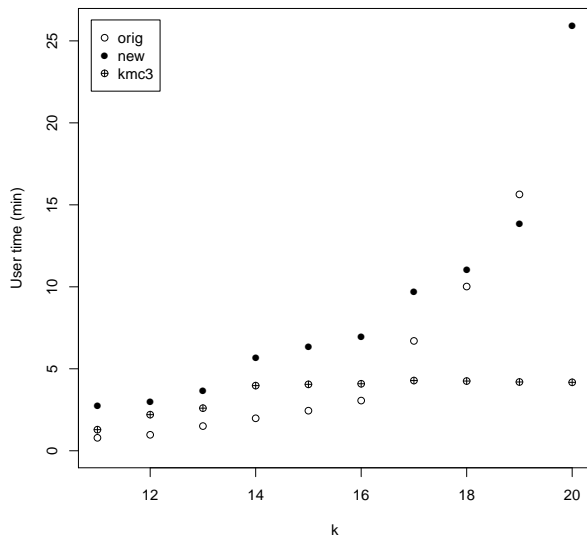


Proof of concept

F. vesca



F. vesca



Summary

- ▶ $2k - 1$ bit encoding of canonical k -mers
- ▶ including palindromes
- ▶ constant time transformation from $2k$ bit to $2k - 1$ bit

Conclusion

Summary

- ▶ $2k - 1$ bit encoding of canonical k -mers
- ▶ including palindromes
- ▶ constant time transformation from $2k$ bit to $2k - 1$ bit

Questions

- ▶ Other encodings?
- ▶ Useful?

Conclusion

Summary

- ▶ $2k - 1$ bit encoding of canonical k -mers
- ▶ including palindromes
- ▶ constant time transformation from $2k$ bit to $2k - 1$ bit

Questions

- ▶ Other encodings?
- ▶ Useful?

Thank you very much!

Appendix

New encoding (k odd)

Encode from the outside inwards

A G A T G A T
11 10 1 11

unspecific pair

A ... T \mapsto 11 ... 11

C ... G \mapsto 11 ... 10

G ... C \mapsto 11 ... 01

T ... A \mapsto 11 ... 00

specifying case

A ... A \mapsto 00 ... 0

A ... C \mapsto 00 ... 1

A ... G \mapsto 01 ... 0

C ... A \mapsto 01 ... 1

C ... C \mapsto 10 ... 0

G ... A \mapsto 10 ... 1

A \mapsto 0

C \mapsto 1

New encoding (k odd)

Encode from the outside inwards

A G A T G A T
11 10 00 11 10 1 11

unspecific pair

A ... T \mapsto 11 ... 11

C ... G \mapsto 11 ... 10

G ... C \mapsto 11 ... 01

T ... A \mapsto 11 ... 00

specifying case

A ... A \mapsto 00 ... 0

A ... C \mapsto 00 ... 1

A ... G \mapsto 01 ... 0

C ... A \mapsto 01 ... 1

C ... C \mapsto 10 ... 0

G ... A \mapsto 10 ... 1

A \mapsto 0

C \mapsto 1

remainder

A \mapsto 00

C \mapsto 01

G \mapsto 10

T \mapsto 11

New encoding (k odd)

Encode from the outside inwards

$$\begin{array}{cccccc} \mathbf{A} & \mathbf{G} & \mathbf{A} & \mathbf{T} & \mathbf{G} & \mathbf{A} & \mathbf{T} & = & \overline{\mathbf{A} & \mathbf{T} & \mathbf{C} & \mathbf{A} & \mathbf{T} & \mathbf{C} & \mathbf{T}} \\ 11 & 10 & 00 & 11 & 10 & 1 & 11 & & 11 & & & & & & 11 \end{array}$$

unspecific pair

$A \cdots T \mapsto 11 \cdots 11$

$C \cdots G \mapsto 11 \cdots 10$

$G \cdots C \mapsto 11 \cdots 01$

$T \cdots A \mapsto 11 \cdots 00$

specifying case

$A \cdots A \mapsto 00 \cdots 0$

$A \cdots C \mapsto 00 \cdots 1$

$A \cdots G \mapsto 01 \cdots 0$

$C \cdots A \mapsto 01 \cdots 1$

$C \cdots C \mapsto 10 \cdots 0$

$G \cdots A \mapsto 10 \cdots 1$

$A \mapsto 0$

$C \mapsto 1$

remainder

$A \mapsto 00$

$C \mapsto 01$

$G \mapsto 10$

$T \mapsto 11$

New encoding (k odd)

Encode from the outside inwards

$$\begin{array}{cccccc} \mathbf{A} & \mathbf{G} & \mathbf{A} & \mathbf{T} & \mathbf{G} & \mathbf{A} & \mathbf{T} & = & \overline{\mathbf{A} & \mathbf{T} & \mathbf{C} & \mathbf{A} & \mathbf{T} & \mathbf{C} & \mathbf{T}} \\ 11 & 10 & 00 & 11 & 10 & 1 & 11 & & 11 & 10 & & & 1 & 11 \end{array}$$

unspecific pair

A ... T \mapsto 11 ... 11

C ... G \mapsto 11 ... 10

G ... C \mapsto 11 ... 01

T ... A \mapsto 11 ... 00

specifying case

A ... A \mapsto 00 ... 0

T ... T \mapsto 00 ... 1

A ... C \mapsto 00 ... 1

G ... T \mapsto 01 ... 0

A ... G \mapsto 01 ... 0

C ... T \mapsto 01 ... 1

C ... A \mapsto 01 ... 1

T ... G \mapsto 10 ... 0

C ... C \mapsto 10 ... 0

G ... G \mapsto 10 ... 1

G ... A \mapsto 10 ... 1

T ... C \mapsto 10 ... 1

A \mapsto 0

T \mapsto 0

C \mapsto 1

G \mapsto 1

remainder

A \mapsto 00

C \mapsto 01

G \mapsto 10

T \mapsto 11

New encoding (k odd)

Encode from the outside inwards

$$\begin{array}{cccccc} \mathbf{A} & \mathbf{G} & \mathbf{A} & \mathbf{T} & \mathbf{G} & \mathbf{A} & \mathbf{T} & = & \overline{\mathbf{A} & \mathbf{T} & \mathbf{C} & \mathbf{A} & \mathbf{T} & \mathbf{C} & \mathbf{T}} \\ 11 & 10 & 00 & 11 & 10 & 1 & 11 & & 11 & 10 & 00 & 11 & 10 & 1 & 11 \end{array}$$

unspecific pair

A ... T \mapsto 11 ... 11

C ... G \mapsto 11 ... 10

G ... C \mapsto 11 ... 01

T ... A \mapsto 11 ... 00

specifying case

A ... A \mapsto 00 ... 0

T ... T \mapsto 00 ... 1

A ... C \mapsto 00 ... 1

G ... T \mapsto 01 ... 0

A ... G \mapsto 01 ... 0

C ... T \mapsto 01 ... 1

C ... A \mapsto 01 ... 1

T ... G \mapsto 10 ... 0

C ... C \mapsto 10 ... 0

G ... G \mapsto 10 ... 1

G ... A \mapsto 10 ... 1

T ... C \mapsto 11 ... 0

A \mapsto 0

T \mapsto 0

C \mapsto 1

G \mapsto 1

reverse complement

remainder

A \mapsto 00

C \mapsto 01

G \mapsto 10

T \mapsto 11

New encoding (k odd)

Encode from the outside inwards

$$\begin{array}{cccccc} \mathbf{A} & \mathbf{G} & \mathbf{A} & \mathbf{T} & \mathbf{G} & \mathbf{A} & \mathbf{T} & = & \overline{\mathbf{A} & \mathbf{T} & \mathbf{C} & \mathbf{A} & \mathbf{T} & \mathbf{C} & \mathbf{T}} \\ 11 & 10 & 00 & 11 & 10 & 1 & 11 & & 11 & 10 & 00 & 11 & 10 & 1 & 11 \end{array}$$

unspecific pair

A ... T \mapsto 11 ... 11

C ... G \mapsto 11 ... 10

G ... C \mapsto 11 ... 01

T ... A \mapsto 11 ... 00

specifying case

A ... A \mapsto 00 ... 0

T ... T \mapsto 00 ... 1

A ... C \mapsto 00 ... 1

G ... T \mapsto 01 ... 0

A ... G \mapsto 01 ... 0

C ... T \mapsto 01 ... 1

C ... A \mapsto 01 ... 1

T ... G \mapsto 10 ... 0

C ... C \mapsto 10 ... 0

G ... G \mapsto 10 ... 1

G ... A \mapsto 10 ... 1

T ... C \mapsto 10 ... 1

A \mapsto 0

T \mapsto 0

C \mapsto 1

G \mapsto 1

reverse complement

remainder

A \mapsto 00

C \mapsto 01

G \mapsto 10

T \mapsto 11

Bijectiveness

► decoding

New encoding (k odd)

00000	AAA	10000	CAC
00001	AAC	10001	GAA
00010	ACA	10010	CCC
00011	ACC	10011	GCA
00100	AGA	10100	CGC
00101	AGC	10101	GGA
00110	ATA	10110	CTC
00111	ATC	10111	GTA
01000	AAG	11000	TAA
01001	CAA	11001	GAC
01010	ACG	11010	CAG
01011	CCA	11011	AAT
01100	AGG	11100	TCA
01101	CGA	11101	GCC
01110	ATG	11110	CCG
01111	CTA	11111	ACT

Efficient implementation

2-bit encoding: shift/update (const. time)

A G A T A T \mapsto 00 10 00 11 00 11

G A T A T A \mapsto 10 00 11 00 11 00

Efficient implementation

2-bit encoding: shift/update (const. time)

AGATAT \mapsto 00 10 00 11 00 11
GATATA \mapsto 10 00 11 00 11 00

new encoding: shift/update (const. time) + shorten (const. time)

AGATAT \mapsto 00 10 00 11 00 11 \mapsto 11 10 00 11 1 11
GATATA \mapsto 10 00 11 00 11 01 \mapsto 10 00 11 00 11 1

Efficient implementation

Shorten: 2-bit \Rightarrow new encoding (const. time)

k -mer x	A C C C T C G T
\bar{x}	A C G A G G G T
2-bit: x	0001010111011011
\bar{x}	0001100010101011

_111110011101011

Efficient implementation

Shorten: 2-bit \Rightarrow new encoding (const. time)

k -mer x A C C C T C G T

\bar{x} A C G A G G G T

2-bit: x 0001010111011011

\bar{x} 0001100010101011

XOR 0000110101110000

count leading zeros 4

111110011101011

Efficient implementation

Shorten: 2-bit \Rightarrow new encoding (const. time)

k -mer x A C C C T C G T

\bar{x} A C G A G G G T

2-bit: x 0001010111011011

\bar{x} 0001100010101011

XOR 0000110101110000

count leading zeros 4

0000000000001011

-11110101110000

OR -111101011101011

set specifying pos. -11110011101011

Even-length k -mers

Encoding palindromes

A G A T C T
11 11 01 11

unspecific pair

A ... T \mapsto 11 ... 11

C ... G \mapsto 11 ... 10

G ... C \mapsto 11 ... 01

T ... A \mapsto 11 ... 00

specifying case

A ... A \mapsto 00 ... 0

A ... C \mapsto 00 ... 1

A ... G \mapsto 01 ... 0

C ... A \mapsto 01 ... 1

C ... C \mapsto 10 ... 0

G ... A \mapsto 10 ... 1

A \mapsto 0

C \mapsto 1

remainder

A \mapsto 00

C \mapsto 01

G \mapsto 10

T \mapsto 11

Even-length k -mers

Encoding palindromes

A G A T C T
11 11 01 11

unspecific pair

A ... T \mapsto 11 ... 11

C ... G \mapsto 11 ... 10

G ... C \mapsto 11 ... 01

T ... A \mapsto 11 ... 00

specifying case

A ... A \mapsto 00 ... 0

A ... C \mapsto 00 ... 1

A ... G \mapsto 01 ... 0

C ... A \mapsto 01 ... 1

C ... C \mapsto 10 ... 0

G ... A \mapsto 10 ... 1

A \mapsto 0

C \mapsto 1

palindrome identified

AT \mapsto 011 ... 111

CG \mapsto 011 ... 110

GC \mapsto 100 ... 001

TA \mapsto 100 ... 000

Even-length k -mers

Encoding palindromes

A G A T C T
0 11 11 111 01 11

unspecific pair

A ... T \mapsto 11 ... 11

C ... G \mapsto 11 ... 10

G ... C \mapsto 11 ... 01

T ... A \mapsto 11 ... 00

specifying case

A ... A \mapsto 00 ... 0

A ... C \mapsto 00 ... 1

A ... G \mapsto 01 ... 0

C ... A \mapsto 01 ... 1

C ... C \mapsto 10 ... 0

G ... A \mapsto 10 ... 1

A \mapsto 0

C \mapsto 1

palindrome identified

AT \mapsto 011 ... 111

CG \mapsto 011 ... 110

GC \mapsto 100 ... 001

TA \mapsto 100 ... 000

Even-length k -mers

Encoding palindromes

A G A T C T
0 11 11 111 01 11

A G G C C T
1 00 00 001 01 11

unspecific pair

A ... T \mapsto 11 ... 11

C ... G \mapsto 11 ... 10

G ... C \mapsto 11 ... 01

T ... A \mapsto 11 ... 00

specifying case

A ... A \mapsto 00 ... 0

A ... C \mapsto 00 ... 1

A ... G \mapsto 01 ... 0

C ... A \mapsto 01 ... 1

C ... C \mapsto 10 ... 0

G ... A \mapsto 10 ... 1

A \mapsto 0

C \mapsto 1

palindrome identified

AT \mapsto 011 ... 111

CG \mapsto 011 ... 110

GC \mapsto 100 ... 001

TA \mapsto 100 ... 000