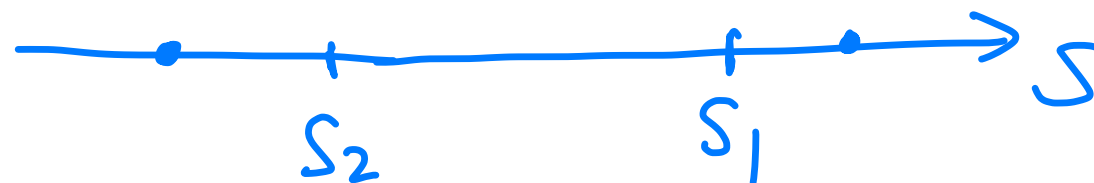# CSE 566 Spring 2023

## Minimizers

Instructor: Mingfu Shao

# Gapped LSH Family

- A set of hash functions $\mathcal{F}$ is said to be a $(s_1, s_2, p_1, p_2)$-*gapped LSH family* for similarity measure $s(\cdot, \cdot)$, where $s_1 \geq s_2$ and $p_1 \geq p_2$, if for any two $x$ and $y$ we have

  - If $s(x, y) \geq s_1$, then $\Pr_{f \in \mathcal{F}}(f(x) = f(y)) \geq p_1$;

  - If $s(x, y) \leq s_2$, then $\Pr_{f \in \mathcal{F}}(f(x) = f(y)) \leq p_2$.

- If $\mathcal{F}$ is a $(r, r, r, r)$-*gapped LSH family for measure $s(\cdot, \cdot)$* for every $0 < r < 1$, then $\mathcal{F}$ is a LSH family for measure $s(\cdot, \cdot)$

# Results for Edit Distance

X = ACGTGTAC
Y = ACTGTAC

- Define $es(x, y) := 1 - d(x, y)/n$ as the edit-similarity between strings $x$ and $y$ of length $n$, where $d(x, y)$ is the edit distance.

- Edit distance/similarity is fundamentally different from the distance/similarity on normed vector space.

X = 0 1 1 0
Y = 1 0 1 1

X = {1, 0, ... 1}
Y = {0, 1, ... 1}
U = {x_1, ... x_n}

X = ( ) ∈ ℝ^d.
Y = 1 ( ) ∈ ℝ^d

- It remains open whether a LSH family exists for edit-similarity.

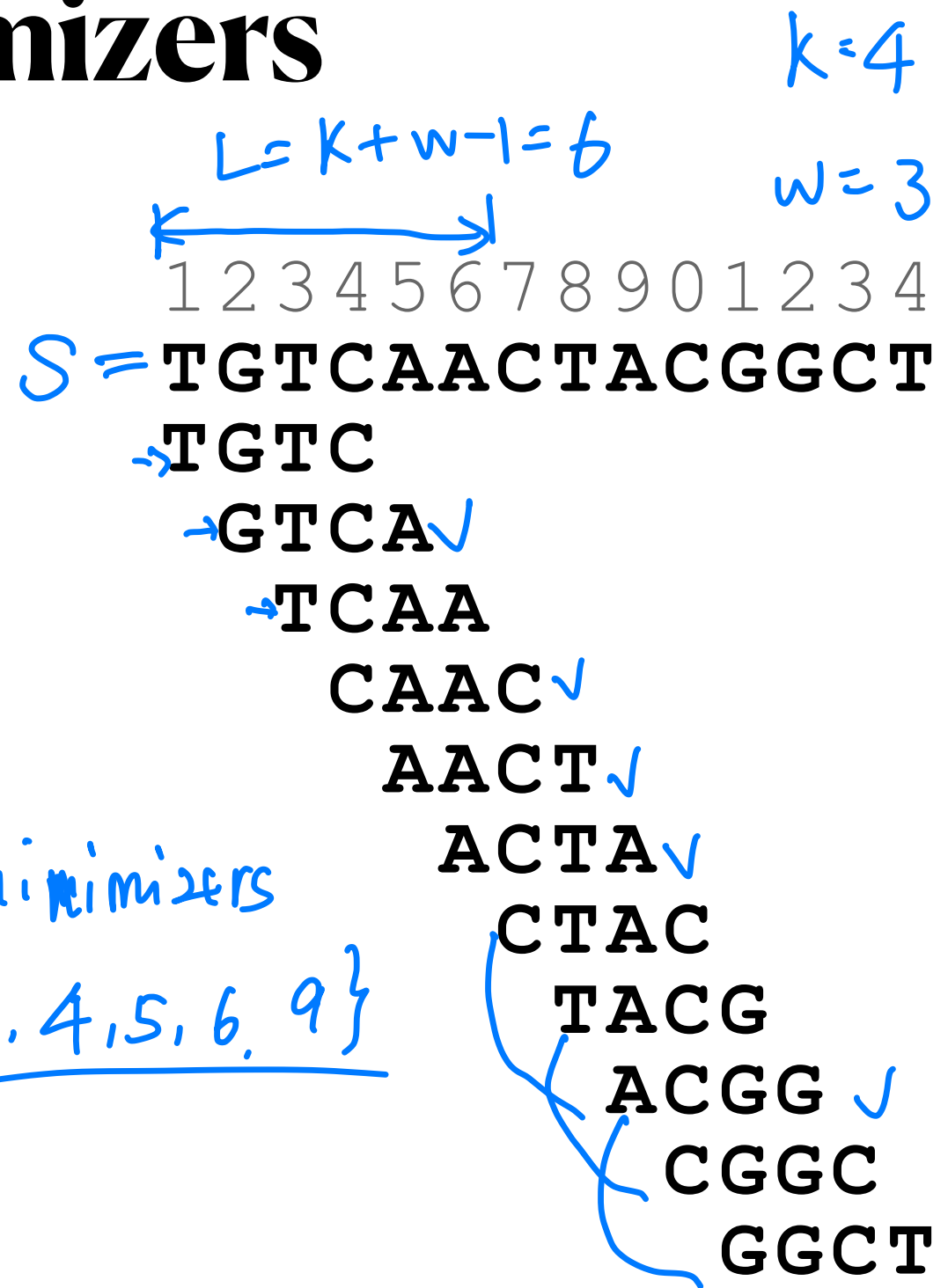- A gapped LSH family for the edit distance exists!

  - OMH: https://doi.org/10.1093/bioinformatics/btz354 (2019)

order min hash

# Minimizers

$k = 4$

$w = 3$

$L = k + w - 1 = 6$

- Given a string $S$, the Minimizers select the smallest kmer (according to order $\pi$) in each sliding window of $w$ kmers.

  - $k$: the length of kmer
  - $w$: the window size
  - $\pi$: an order over all kmers

AAAA, AAAC, ..., TTTT

```
        1 2 3 4 5 6 7 8 9 0 1 2 3 4
S =     T G T C A A C T A C G G C T
→       T G T C
→       G T C A ✓
→       T C A A
        C A A C ✓
        A A C T ✓
        A C T A ✓
        C T A C
        T A C G
        A C G G ✓
        C G G C
        G G C T
```

minimizers

$\{2, 4, 5, 6, 9\}$

# Calculating Minimizers

```
Input: string S, k, w, order pi

Output: array M to store the positions of minimizers

init an empty queue Q

for i = 1 to (|S| - k + 1)

    let m = S[i..i+k-1] be the current kmer

    while m <(pi) S[tail(Q)..tail(Q)-k+1]: pop-tail(Q)

    append i to the tail of Q

    while head(Q) <= i - w: pop-head(Q)

    if i >= w && head(Q)!= tail(M): append head(Q) to M

end for

return M;
```
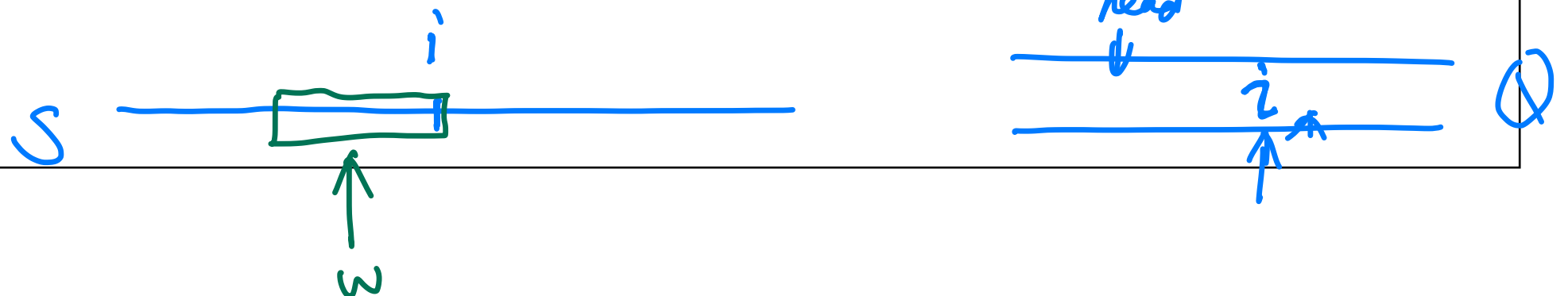
kmer at the tail of Q

S

i

w

head

i

Q

# An Example

↓ ↓↓ ↓

1234567890 1234

**TGTCAACTACGGCT**
**TGTC**
  **GTCA**
    **TCAA**
      **CAAC**
        **AACT**
          **ACTA**
            **CTAC**
              **TACG**
                **ACGG**
                  **CGGC**
                    **GGCT**

Q: $\sharp 4$
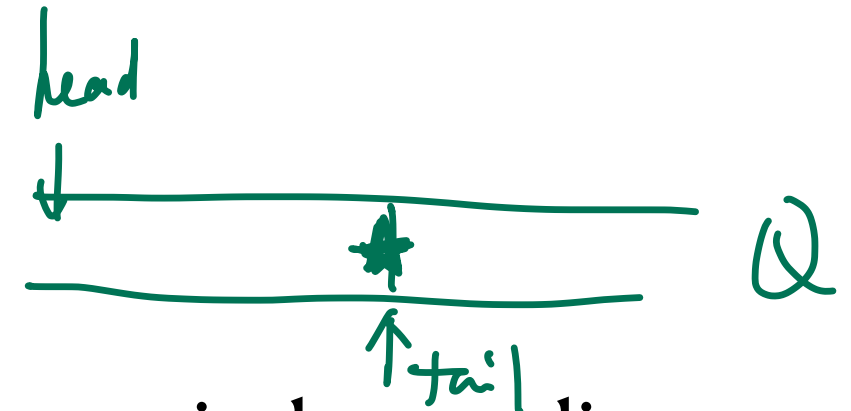
M $[2, 4 ..$

# Correctness

- Claim: head(Q) is the minimizer of the current window ending with the i-th kmer
  - Fact 1: it is safe to pop-tail(Q) when m is smaller than the kmer at tail(Q)
  - Fact 2: it is safe to pop-head(Q) when it is out of the current window
  - Fact 3: the kmers in Q are in increasing order
  - Fact 4: head(Q) is the smallest kmer in Q

# Running Time

- The entire algorithm takes $O(|S|)$ comparisons

  - Each kmer gets added into $Q$ once

  - Hence #pops is also bounded by $|S|$

$$\#pops \leq \#added\text{-}kmers$$

# Properties of Minimizers

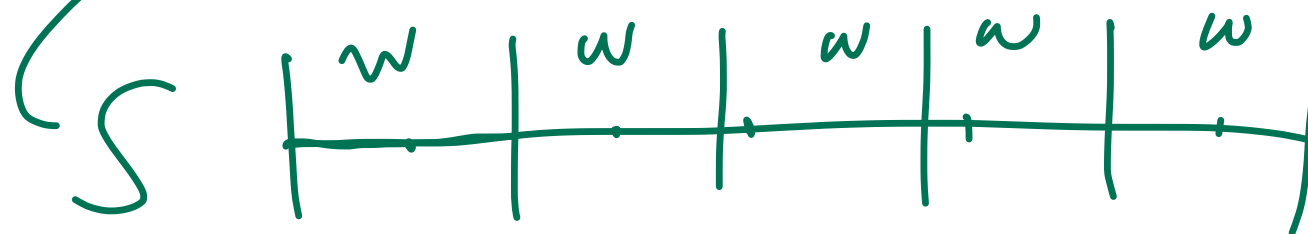Handwritten annotations: $L = w+k-1$, $S_1$, $S_2$, $L$, $w+k-1$

- Small gap: the distance between two selected positions is at most $w$, and hence provides a good "coverage".

- Consistency: if two sequences share a substring of length $w + k - 1$, then they share a minimizer.

- Locality-sensitive: if two sequences share a "similar" substring of length $w + k - 1$, then with high probability that they also share a minimizer.

# Density of Minimizers

- Given $w$, $k$, and and order $\pi$, the <span style="color:blue">particular density</span> of a sequence $S$ is defined as: $|M(S)|/(|S| - k + 1)$.

- Given $w$, $k$, and and order $\pi$, the <span style="color:blue">expected density</span> is defined as the particular density over an infinity-length, random sequence (i.e., each base is sampled uniformly at random).

- Trivial bounds for expected density: $1/w \leq d \leq 1$.

- Proved lower bound for expected density: ~~$1/(w+1) \leq d$~~   $d \geq \dfrac{1.5}{1+w}$

# Density with Random Order

- Order $\pi$ is picked uniformed at random from all orderings.

- **Theorem**: the expected density wrt a random order $\pi$ is $2/(1 + w) + o(1/w)$.

- **Proof**: consider the probability that the next window uses a new minimizer.

$$Pr = \frac{2}{1+w}$$