

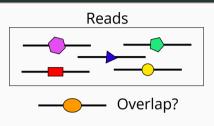


# Locality sensitive hashing for the edit distance

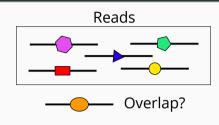
Guillaume Marçais, Dan DeBlasio, Prashant Pandey, Carl Kingsford

Carnegie Mellon University

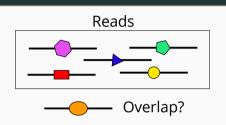
- Compute overlaps between reads (HMAP)
- Instance of "Nearest Neighbor Problem" for edit distance
- Use multiple hash tables
- Need meaningful hash collisions



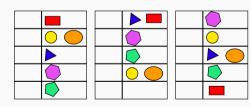
- Compute overlaps between reads (HMAP)
- Instance of "Nearest Neighbor Problem" for edit distance
- Use multiple hash tables
- Need meaningful hash collisions



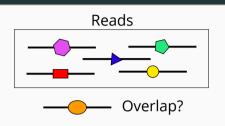
- Compute overlaps between reads (HMAP)
- Instance of "Nearest Neighbor Problem" for edit distance
- Use multiple hash tables
- Need meaningful hash collisions



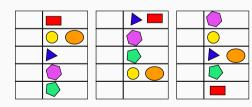
#### **Hash Tables**



- Compute overlaps between reads (HMAP)
- Instance of "Nearest Neighbor Problem" for edit distance
- Use multiple hash tables
- Need meaningful hash collisions



#### **Hash Tables**



# **Locality Sensitive Hashing**

Pick h at random from  $\mathcal{H}$ :

$$\Pr[h(\bigcirc) = h(\bigcirc)] > \Pr[h(\bigcirc) = h(\bigcirc)]$$

## **Locality sensitive hash family**

Family  ${\cal H}$  of hash functions where similar elements are more likely to have the same value than distant elements.

# **Locality Sensitive Hashing**

Pick h at random from  $\mathcal{H}$ :

$$\Pr[h(\bigcirc) = h(\bigcirc)] > \Pr[h(\bigcirc) = h(\triangleright)]$$

$$Sketch(\bigcirc) = \{h_1(\bigcirc), \ldots, h_m(\bigcirc)\}$$

## **Locality sensitive hash family**

Family  ${\cal H}$  of hash functions where similar elements are more likely to have the same value than distant elements.

# **Locality Sensitive Hashing**

The family  $\mathcal{H}$  is sensitive for distance  $\mathrm{D}$  if there exists  $d_1 < d_2$ ,  $p_1 > p_2$  such that for all  $x,y \in \mathcal{U}$ 

$$D(x, y) \le d_1 \implies \Pr_{h \in \mathcal{H}}[h(x) = h(y)] \ge p_1$$
  
 $D(x, y) \ge d_2 \implies \Pr_{h \in \mathcal{H}}[h(x) = h(y)] \le p_2$ 

- Low distance ←⇒ High collisions
- ullet High distance  $\Longleftrightarrow$  Low collisions

#### Locality sensitive hash family

Family  $\mathcal{H}$  of hash functions where similar elements are more likely to have the same value than distant elements.

## LSH for the edit distance

## How to design an LSH for edit distance?

- LSH for Jaccard distance (minHash) used as proxy
- Jaccard distance is **significantly different** than edit distance

#### LSH for the edit distance

## How to design an LSH for edit distance?

- LSH for Jaccard distance (minHash) used as proxy
- Jaccard distance is **significantly different** than edit distance

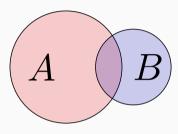
## LSH for the edit distance

## How to design an LSH for edit distance?

- LSH for Jaccard distance (minHash) used as proxy
- Jaccard distance is **significantly different** than edit distance

# **Jaccard distance**

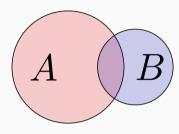
## Jaccard distance between sets A,B:



$$J(A, B) = 1 - \frac{|A \cap B|}{|A \cup B|}$$

## **Jaccard distance**

## Jaccard distance between sets A, B:



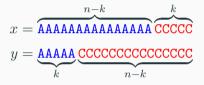
$$J(A,B) = 1 - \frac{|A \cap B|}{|A \cup B|}$$

Jaccard between sequences x, y: laccard distance of their k-mer sets

$$J(x,y) = J(\mathcal{K}(x), \mathcal{K}(y))$$

- Low  $D(x,y) \implies \text{Low } J(x,y)$
- $\bullet \ \ \mathsf{High} \ \mathrm{D}(x,y) \ \Longrightarrow \ \ \mathsf{High} \ \mathrm{J}(x,y)$

# Jaccard ignores k-mer repetition



# Jaccard ignores k-mer repetition

# Jaccard ignores *k*-mer repetition

Jaccard distance 
$$\mathrm{J}(x,y)=0$$
 Edit distance  $\mathrm{D}(x,y)\geq 1-\frac{2k}{n}$  Identical  $k$ -mer content and high edit distance

# Weighted Jaccard handles repetitions

Weighted Jaccard 
$${
m J^w}(x,y)=1-\frac{k+2}{n}$$
 Edit distance  ${
m D}(x,y)\geq 1-\frac{2k}{n}$  Weighted Jaccard = Jaccard for multi-sets

# Jaccard and weighted Jaccard ignore relative order

```
x = \mathtt{CCCCACCAACACAAAACCC}
```

 $y = \mathtt{AAAACACAACCCCACCAAA}$ 

# Jaccard and weighted Jaccard ignore relative order

```
x = \texttt{CCCCACCAACAACACAAACCC} \\ \rightarrow \left\{ \substack{\text{AAAA,AAAC,AACA,AACC,ACAA,ACCC,ACCA,CCCC},\\ \text{CAAA,CAAC,CACA,CACC,CCCA,CCCC}} \right\} \\ y = \texttt{AAAACACAACCCCACCAAAA} \\ \rightarrow \left\{ \substack{\text{AAAA,AAAC,AACA,AACC,ACAA,ACCC,ACCA,ACCC},\\ \text{CAAA,CAAC,CACA,CACC,CCCA,CCCC}} \right\}
```

x,y: de Bruijn sequences, contain all 16 possible 4-mers once

# Jaccard and weighted Jaccard ignore relative order

$$x = \texttt{CCCCACCAACACACAAAACCC}$$
 $y = \texttt{AAAACACAACCCCACCAAA}$ 

- $\rightarrow \left\{ \begin{smallmatrix} AAAA,AAAC,AACA,AACC,ACAA,ACAC,ACCA,ACCC\\ CAAA,CAAC,CACA,CACC,CCAA,CCAC,CCCA,CCCC \end{smallmatrix} \right\}$
- $\rightarrow \left\{ \begin{matrix} \texttt{AAAA}, \texttt{AAAC}, \texttt{AACA}, \texttt{AACC}, \texttt{ACCA}, \texttt{ACCC}, \\ \texttt{CAAA}, \texttt{CAAC}, \texttt{CACA}, \texttt{CACC}, \texttt{CCCA}, \texttt{CCCC}, \end{matrix} \right\}$

x, y: de Bruijn sequences, contain all 16 possible 4-mers once

$$J(x,y) = J^{w}(x,y) = 0$$
  $D(x,y) = 0.63$ 

## Jaccard is different from edit distance

Unlike edit distance, Jaccard is insensitive to:

- 1. *k*-mer repetitions
- 2. relative positions of k-mers

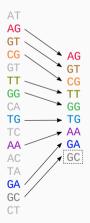
#### **OMH: Order Min Hash**

- minHash is an LSH for Jaccard
- OMH is a refinement of minHash
- OMH is sensitive to
  - ullet repeated k-mers
  - relative order of *k*-mers

 $x = \mathtt{AGTTGAGCGGAAGGTG}, \, k = 2$ 

 $x = \operatorname{AGTTGAGCGGAAGGTG}$ , k = 2

Order: permutation of  $\Sigma^k$ 



x = AGTTGAGCGGAAGGTG, k = 2, m = 6

```
1 2 3 4 5 6
AG GG CG AA TG TT
GT GA GA AG TT GG
CG CG TG TT GG AG
TT AG AG GT CG GC
GG GC GC TG AA GA
TG GT GG GC GT AA
AA AA TT GA AG CG
GA TT AA CG GC TG
GC TG GT GG GA GT
```

$$x = \text{AGTTGAGCGGAAGGTG}, k = 2, m = 6$$

```
1 2 3 4 5 6
AG GG CG AA TG TT
GT GA GA AG TT GG
CG CG TG TT GG AG
TT AG AG GT CG GG
GG GC GC TG AA GA
TG GT GG GC GT AA
AA AA TT GA AG CG
GA TT AA CG GC TG
GC TG GT GG GA GT
```

$$x={\tt AGTTGAGCGGAAGGTG},\, k=2,\, m=6$$
 Order: permutation of  $\Sigma^k\times\{1,\dots,n\}$ 

```
1 2 3 4 5 6
AG GG CG AA TG TT
GT GA GA AG TT GG
CG CG TG TT GG AG
TT AG AG GT CG GC
GG GC GC TG AA GA
TG GT GG GC GT AA
AA AA TT GA AG CG
GA TT AA CG GC TG
GC TG GT GG GA GT
```

```
GA, 4
TG, 3
AG, 5
GT, 1
GT, 13
AA, 10
AG, 11
TT, 2
AG, o
CG, 7
GG, 12
GC, 6
TG, 14
GG, 8
GA, 9
```

$$x = \text{AGTTGAGCGGAAGGTG}, k = 2, m = 6$$

```
1 2 3 4 5 6
AG GG CG AA TG TT
GT GA GA AG TT GG
CG CG TG TT GG AG
TT AG AG GT CG GC
GG GC GC TG AA GA
TG GT GG GC GT AA
AA AA TT GA AG CG
GA TT AA CG GC TG
GC TG GT GG GA GT
```

```
GA, 4 CG, 7 GT, 13 AG, 0 AA, 10
     TG, 14 GA, 4 TT, 2 GT, 13
     AG, 0 GA, 9 AG, 11 GA, 9
GT, 1 GA, 9 TG, 3 AG, 5 GT, 1 TG, 14
GT, 13 AG, 5 AG, 5 AA, 10 AG, 5 GT, 13
     AG, 11 CG, 7 GT, 13 TT, 2 TT, 2
AA, 10
AG, 11 GA, 4 TT, 2 CG, 7 GA, 4 AA, 10
TT, 2 GT, 13 AA, 10 GG, 8 CG, 7 AG, 0
AG, 0 TT, 2 GG, 12 GA, 4 AG, 0 CG, 7
CG, 7 TG, 3 GG, 8 GA, 9 TG, 3 GG, 12
GG, 12 GG, 8 TG, 14 TG, 14 GG, 8 AG, 11
GC, 6 AA, 10 GT, 1 TG, 3 GG, 12 TG, 3
TG, 14 GG, 12 AG, 11 GC, 6 GC, 6 GT, 1
GG, 8 GT, 1 GC, 6 GT, 1 AG, 11 GA, 4
GA, 9 GC, 6 AG, 0 GG, 12 TG, 14 AG, 5
```

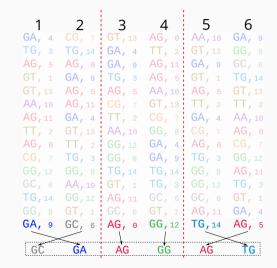
$$x = \text{AGTTGAGCGGAAGGTG}, k = 2, m = 6$$

1 2 3 4 5 6
AG GG CG AA TG TT
GT GA GA AG TT GG
CG CG TG TT GG AG
TT AG AG GT CG GC
GG GC GC TG AA GA
TG GT GG GC GT AA
AA AA TT GA AG CG
GA TT AA CG GC TG
GC TG GT GG GA GT

```
3 4
GA, 4 CG, 7 GT, 13 AG, 0 AA, 10 GA, 9
TG, 3 TG, 14 GA, 4 TT, 2 GT, 13 GG, 8
AG, 5 AG, 0 GA, 9 AG, 11 GA, 9 GC, 6
GT, 1 GA, 9 TG, 3 AG, 5 GT, 1 TG, 14
AG, 11 GA, 4 TT, 2 CG, 7 GA, 4 AA, 10
AG, 0 TT, 2 GG, 12 GA, 4 AG, 0 CG, 7
CG, 7 TG, 3 GG, 8 GA, 9 TG, 3 GG, 12
GG, 8 GT, 1 GC, 6 GT, 1 AG, 11 GA, 4
GA, 9 GC, 6 AG, 0 GG, 12 TG, 14 AG, 5
```

$$x = {\tt AGTTGAGCGGAAGGTG}, \, k = 2, \, m = 6, \, \ell = 2$$

1	2	3	4	5	6
			AA		
	GA	GA			
				AA	GA
					AA
AA	AA		GA		
GA		AA			
GC	TG	GT	GG	GA	GT



#### **OMH** is a LSH for edit distance

## Theorem: OMH is a LSH for edit distance

There exists  $(d_1, d_2, p_1, p_2)$  such that OMH is sensitive for the edit distance.

## Conclusion

- OMH:
  - improvement on minHash
  - easy to compute
  - locality sensitive for edit distance
- LSH for other alignment scores?
- Smallest "gap" achievable?

## Conclusion

- OMH:
  - improvement on minHash
  - easy to compute
  - locality sensitive for edit distance
- LSH for other alignment scores?
- Smallest "gap" achievable?

## Conclusion

- OMH:
  - improvement on minHash
  - easy to compute
  - locality sensitive for edit distance
- LSH for other alignment scores?
- Smallest "gap" achievable?



# Thank you

Laura Tung Hongyu Zheng Yutong Qiu

> Minh Hoang Mohsen Ferdosi

Yihang Shen

Natalie Sauerwald

Hongyu Zheng

Cong Ma

Shawn Baker Yinjie Gao







R01GM122935

The Shurl and Kay Curci Foundation