



Sketching and locality sensitive hashing for alignment

Guillaume Marçais

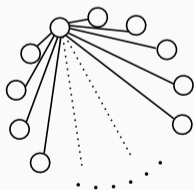
2/15/23

Why do we need sketching and Locality Sensitive Hashing for alignment?

Large scale alignment problems

Cluster N samples based on sequence similarity

- $\rightarrow N^2/2$ alignment problems
- Speed-up pairwise alignment task?
- Skip hopeless alignments?



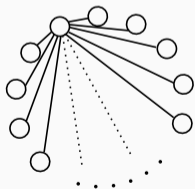
Sequence search in large database

- Avoid aligning to all sequences in database?
- Approximate nearer neighbor search
- High dimension, non-geometric space

Large scale alignment problems

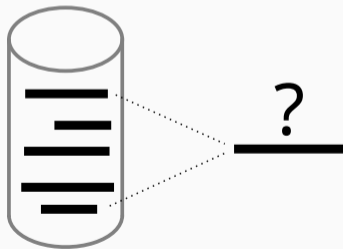
Cluster N samples based on sequence similarity

- $\rightarrow N^2/2$ alignment problems
- Speed-up pairwise alignment task?
- Skip hopeless alignments?

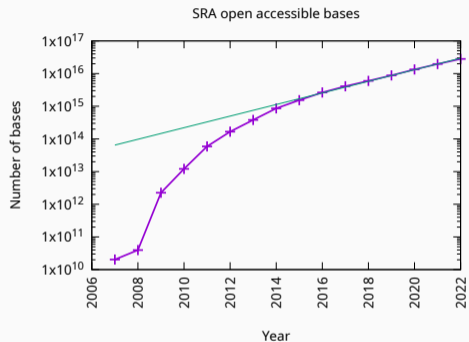


Sequence search in large database

- Avoid aligning to all sequences in database?
- Approximate nearer neighbor search
- High dimension, non-geometric space



Fast growth of sequence databases



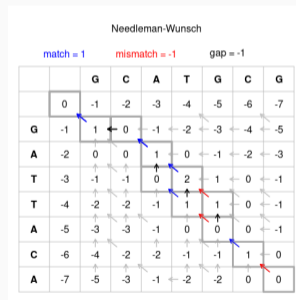
- Exponential growth in public and private databases (SRA: $1.5 \times$ /year)
- \implies hidden exponential slow down in large scale analysis

Sequence alignment is hard

No strongly subquadratic time algorithm, most likely (Backurs, Indyk 2015)

Computing the edit distance E_d in time $O(n^{2-\delta})$, $\delta > 0$ violates the Strong Exponential Time Hypothesis (SETH).

- Usual dynamic programming: $O(n^2)$
- ¹Masek and Paterson: $O\left(\frac{n^2}{\log(n)}\right)$
- $n^{2-\delta} \ll \frac{n^2}{\log(n)} \ll n^2$
- Can't fundamentally improve



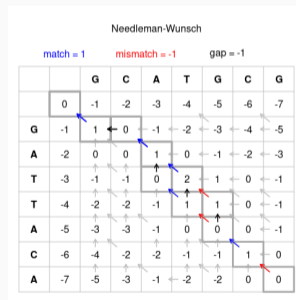
¹A faster algorithm computing string edit distances (1980)

Sequence alignment is hard

No strongly subquadratic time algorithm, most likely (Backurs, Indyk 2015)

Computing the edit distance E_d in time $O(n^{2-\delta})$, $\delta > 0$ violates the Strong Exponential Time Hypothesis (SETH).

- Usual dynamic programming: $O(n^2)$
- ¹Masek and Paterson: $O\left(\frac{n^2}{\log(n)}\right)$
- $n^{2-\delta} \ll \frac{n^2}{\log(n)} \ll n^2$
- Can't fundamentally improve



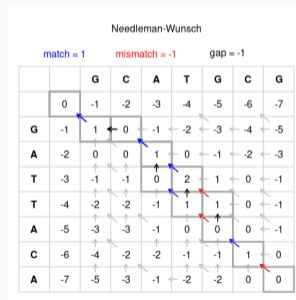
¹A faster algorithm computing string edit distances (1980)

Sequence alignment is hard

No strongly subquadratic time algorithm, most likely (Backurs, Indyk 2015)

Computing the edit distance E_d in time $O(n^{2-\delta})$, $\delta > 0$ violates the Strong Exponential Time Hypothesis (SETH).

- Usual dynamic programming: $O(n^2)$
- ¹Masek and Paterson: $O\left(\frac{n^2}{\log(n)}\right)$
- $n^{2-\delta} \ll \frac{n^2}{\log(n)} \ll n^2$
- Can't fundamentally improve

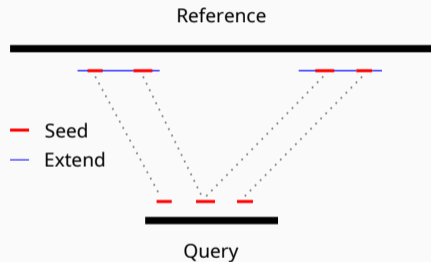


¹A faster algorithm computing string edit distances (1980)

Seed and extend paradigm

Main paradigm:

- Find seeds (small exact matches)
- Cluster “coherent” seeds
- Extend between seeds using DP
- Used since the 90s’ (Blast, MUMmer)
- Still computationally intensive for large scale
- Many ways to find seeds:
 - k -mers
 - Suffix trees/arrays, FM Index
 - LSH / sketching



Avoid computing edit distance directly, use proxy measures easier to compute

- LSH: hashing method to avoid fruitless comparisons
- Sketching: sparse representation allowing quick comparison

Locality Sensitive Hashing: Make collisions matters

\mathcal{U} : universe. T : hash table. $|T| \ll |\mathcal{U}|$. $h : \mathcal{U} \rightarrow [0, |T| - 1]$.

$\mathcal{H} = \{h : \mathcal{U} \rightarrow [0, |T| - 1]\}$

$S = \text{AACGGTG}$

$h(S) = 2$

	T
0	
1	
2	S
3	
4	

Universal Hashing

- Collisions as rare as possible
- $\forall x, y \in \mathcal{U}, x \neq y,$

$$\Pr_{h \in \mathcal{H}} [h(x) = h(y)] = \frac{1}{|T|}$$

Locality Sensitive Hashing

- Collision between similar elements
- $\forall x, y \in \mathcal{U}$

$$E_d(x, y) \leq d_1 \implies \Pr_{h \in \mathcal{H}} [h(x) = h(y)] \geq p_1$$

$$E_d(x, y) \geq d_2 \implies \Pr_{h \in \mathcal{H}} [h(x) = h(y)] \leq p_2$$

Locality Sensitive Hashing: Make collisions matters

\mathcal{U} : universe. T : hash table. $|T| \ll |\mathcal{U}|$. $h : \mathcal{U} \rightarrow [0, |T| - 1]$.

$\mathcal{H} = \{h : \mathcal{U} \rightarrow [0, |T| - 1]\}$

$S = \text{AACGGTG}$

$h(S) = 2$

T	
0	
1	
2	S
3	
4	

Universal Hashing

- Collisions as rare as possible
- $\forall x, y \in \mathcal{U}, x \neq y,$

$$\Pr_{h \in \mathcal{H}} [h(x) = h(y)] = \frac{1}{|T|}$$

Locality Sensitive Hashing

- Collision between similar elements
- $\forall x, y \in \mathcal{U}$

$$E_d(x, y) \leq d_1 \implies \Pr_{h \in \mathcal{H}} [h(x) = h(y)] \geq p_1$$

$$E_d(x, y) \geq d_2 \implies \Pr_{h \in \mathcal{H}} [h(x) = h(y)] \leq p_2$$

Locality Sensitive Hashing: Make collisions matters

\mathcal{U} : universe. T : hash table. $|T| \ll |\mathcal{U}|$. $h : \mathcal{U} \rightarrow [0, |T| - 1]$.

$\mathcal{H} = \{h : \mathcal{U} \rightarrow [0, |T| - 1]\}$

$S = \text{AACGGTG}$

$h(S) = 2$

	T
0	
1	
2	S
3	
4	

Universal Hashing

- Collisions as rare as possible
- $\forall x, y \in \mathcal{U}, x \neq y,$

$$\Pr_{h \in \mathcal{H}} [h(x) = h(y)] = \frac{1}{|T|}$$

Locality Sensitive Hashing

- Collision between similar elements
- $\forall x, y \in \mathcal{U}$

$$E_d(x, y) \leq d_1 \implies \Pr_{h \in \mathcal{H}} [h(x) = h(y)] \geq p_1$$

$$E_d(x, y) \geq d_2 \implies \Pr_{h \in \mathcal{H}} [h(x) = h(y)] \leq p_2$$

Locality Sensitive Hashing Definition

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

$$D(x, y) \leq d_1 \implies \Pr_{h \in \mathcal{H}} [h(x) = h(y)] \geq p_1$$

$$D(x, y) \geq d_2 \implies \Pr_{h \in \mathcal{H}} [h(x) = h(y)] \leq p_2$$

- **Low distance** \iff **High collisions**
- **High distance** \iff **Low collisions**
- **In between d_1, d_2 : No guarantee**

Locality sensitive hash family

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

Locality Sensitive Hashing Definition

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

$$D(x, y) \leq d_1 \implies \Pr_{h \in \mathcal{H}} [h(x) = h(y)] \geq p_1$$

$$D(x, y) \geq d_2 \implies \Pr_{h \in \mathcal{H}} [h(x) = h(y)] \leq p_2$$

- **Probability over choice of $h \in \mathcal{H}$, not over the elements x, y**

Locality sensitive hash family

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

Locality Sensitive Hashing Definition

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

$$D(x, y) \leq d_1 \implies \Pr_{h \in \mathcal{H}} [h(x) = h(y)] \geq p_1$$

$$D(x, y) \geq d_2 \implies \Pr_{h \in \mathcal{H}} [h(x) = h(y)] \leq p_2$$

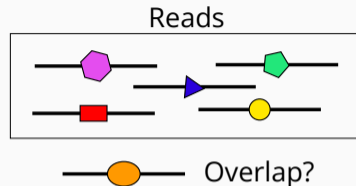
- $d_1 < d_2$: “gapped” LSH
- $d_1 = d_2$, “ungapped” LSH
- Gap not desirable but not always avoidable.

Locality sensitive hash family

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

Overlap computation

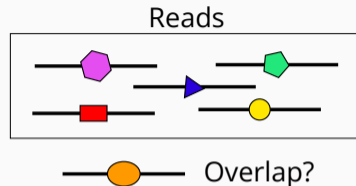
- Compute overlaps between reads (MHAP²)
- Instance of “Nearest Neighbor Problem” for edit distance
- Use multiple hash tables
- Orange ellipse in same location as yellow circle



²Assembling large genomes with single-molecule sequencing and locality-sensitive hashing

Overlap computation

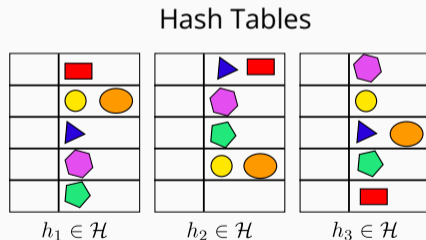
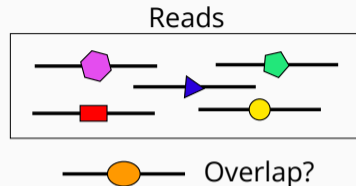
- Compute overlaps between reads (MHAP²)
- Instance of “Nearest Neighbor Problem” for edit distance
- Use multiple hash tables
- Orange ellipse in same location as yellow circle



²Assembling large genomes with single-molecule sequencing and locality-sensitive hashing

Overlap computation

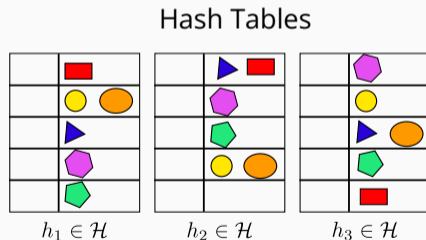
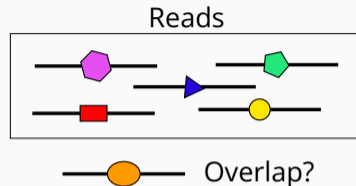
- Compute overlaps between reads (MHAP²)
- Instance of “Nearest Neighbor Problem” for edit distance
- Use multiple hash tables
- Orange ellipse in same location as yellow circle



²Assembling large genomes with single-molecule sequencing and locality-sensitive hashing

Overlap computation

- Compute overlaps between reads (MHAP²)
- Instance of “Nearest Neighbor Problem” for edit distance
- Use multiple hash tables
- Orange ellipse in same location as yellow circle



²Assembling large genomes with single-molecule sequencing and locality-sensitive hashing

How to design an LSH for edit distance?

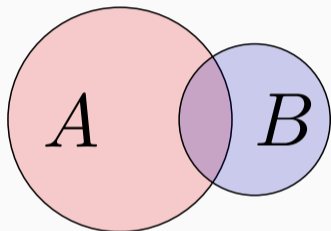
- minHash: LSH for k -mer Jaccard distance
- OMH: Ordered Min Hash

How to design an LSH for edit distance?

- **minHash: LSH for k -mer Jaccard distance**
- **OMH: Ordered Min Hash**

Jaccard distance

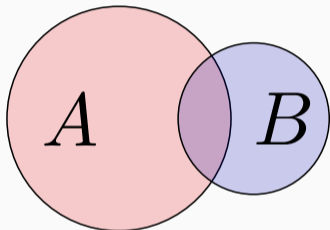
Jaccard distance between sets A, B :



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

Jaccard distance

Jaccard distance between sets A, B :



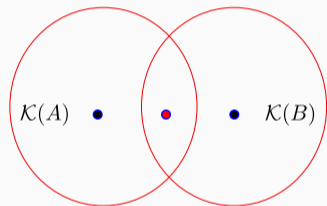
$$J_d(A, B) = 1 - \frac{|A \cap B|}{|A \cup B|}$$

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

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

- Low $E_d(x, y) \implies$ Low $J_d(x, y)$
- High $E_d(x, y) \not\Rightarrow$ High $J_d(x, y)$
- Can have false positive, few false negative

MinHash: an LSH for the Jaccard distance



- Permutation of k -mers: $\pi : 4^k \rightarrow 4^k$ one-to-one

$$\mathcal{H} = \{h_\pi(S) = \operatorname{argmin}_{m \in \mathcal{K}(S)} \pi(m) \mid \pi \text{ permutation of } k\text{-mers}\}$$

- Fix π , every k -mer of $A \cup B$ equally likely to be the minimum for π

$$\Pr_{h \in \mathcal{H}} [h(A) = h(B)] = \frac{|A \cap B|}{|A \cup B|}$$

- Unbiased estimator, ungapped LSH

minHash sketch: dimensionality reduction

- Choose L hash functions from \mathcal{H} : $h_i, 1 \leq i \leq L$
- Sketch of S : vector $\text{Sk}(S) = (h_i(S))_{1 \leq i \leq L}$
- Big compression: Mash³ $L = 1000, k = 21, 7000 \times$ compression
- Very fast pairwise comparison (Hamming distance between sketches)

$$\text{Sk}(A) = \begin{pmatrix} \text{CGAG} \\ \text{TTAC} \\ \text{CATC} \\ \text{CCAT} \\ \text{CATG} \\ \text{ACAA} \end{pmatrix}, \text{Sk}(B) = \begin{pmatrix} \text{GTTT} \\ \text{TTAC} \\ \text{GTAG} \\ \text{ATTT} \\ \text{ACCC} \\ \text{ACAA} \end{pmatrix} \rightarrow J_d(\mathcal{K}(A), \mathcal{K}(B)) \approx 1 - \frac{2}{6}$$

³Mash: fast genome and metagenome distance estimation using MinHash

- minHash: LSH for k -mer Jaccard distance
- **OMH: Ordered Min Hash**

Jaccard ignores k -mer repetition

$$\begin{aligned} x &= \overbrace{\text{AAAAAAAAAAAAAAAA}}^{n-k} \overbrace{\text{CCCC}}^k \\ y &= \overbrace{\text{AAAAA}}^k \overbrace{\text{CCCCCCCCCCCCCCCC}}^{n-k} \end{aligned}$$

Jaccard ignores k -mer repetition

$$\begin{aligned} x &= \overbrace{\text{AAAAAAAAAAAAAAAA}}^{n-k} \overbrace{\text{CCCC}}^k && \rightarrow \{\text{AAAAA}, \text{AAAAC}, \text{AAACC}, \text{AACCC}, \text{ACCCC}, \text{CCCCCC}\} \\ y &= \underbrace{\text{AAAAA}}_k \underbrace{\text{CCCCCCCCCCCCCCC}}_{n-k} && \rightarrow \{\text{AAAAA}, \text{AAAAC}, \text{AAACC}, \text{AACCC}, \text{ACCCC}, \text{CCCCCC}\} \end{aligned}$$

Jaccard ignores k -mer repetition

$$\begin{aligned} x &= \overbrace{\text{AAAAAAAAAAAAAAAA}}^{n-k} \overbrace{\text{CCCC}}^k && \rightarrow \{\text{AAAAA}, \text{AAAAC}, \text{AAACC}, \text{AACCC}, \text{ACCCC}, \text{CCCCCC}\} \\ y &= \underbrace{\text{AAAAA}}_k \underbrace{\text{CCCCCCCCCCCCCCC}}_{n-k} && \rightarrow \{\text{AAAAA}, \text{AAAAC}, \text{AAACC}, \text{AACCC}, \text{ACCCC}, \text{CCCCCC}\} \end{aligned}$$

$$\text{Jaccard distance } J_d(x, y) = 0 \quad \text{Edit distance } E_d(x, y) \geq 1 - \frac{2k}{n}$$

Identical k -mer content and high edit distance

- $\chi_A : \mathcal{U} \rightarrow \{0, 1\}$,
 $\chi_A(x) = 1 \iff x \in A$

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{\sum_{x \in \mathcal{U}} \min(\chi_A(x), \chi_B(x))}{\sum_{x \in \mathcal{U}} \max(\chi_A(x), \chi_B(x))}$$

- $\chi_A^w : \mathcal{U} \rightarrow \mathbb{N}$,
 $\chi_A^w(x) = \#$ of instances of x in A

Weighted Jaccard: Jaccard on multi-set

- $\chi_A : \mathcal{U} \rightarrow \{0, 1\}$,
 $\chi_A(x) = 1 \iff x \in A$

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{\sum_{x \in \mathcal{U}} \min(\chi_A(x), \chi_B(x))}{\sum_{x \in \mathcal{U}} \max(\chi_A(x), \chi_B(x))}$$

- $\chi_A^w : \mathcal{U} \rightarrow \mathbb{N}$,
 $\chi_A^w(x) = \#$ of instances of x in A

$$J^w(A, B) = \frac{\sum_{x \in \mathcal{U}} \min(\chi_A^w(x), \chi_B^w(x))}{\sum_{x \in \mathcal{U}} \max(\chi_A^w(x), \chi_B^w(x))}$$

Weighted Jaccard handles repetitions

$$\begin{aligned}
 x &= \overbrace{\text{AAAAAAAAAAAAAAAA}}^{n-k} \overbrace{\text{CCCC}}^k & \rightarrow & \left\{ \begin{array}{l} (\text{AAAAA},1), (\text{AAAAA},2), \dots, (\text{AAAAA},11) \\ (\text{AAAAC},1), (\text{AAACC},1), (\text{AACCC},1), (\text{ACCCC},1), (\text{CCCCC},1) \end{array} \right\} \\
 y &= \overbrace{\text{AAAAA}}^k \overbrace{\text{CCCCCCCCCCCCCCCC}}^{n-k} & \rightarrow & \left\{ \begin{array}{l} (\text{AAAAA},1), (\text{AAAAC},1), (\text{AAACC},1), (\text{AACCC},1), (\text{ACCCC},1) \\ (\text{CCCCC},1), (\text{CCCCC},2), \dots, (\text{CCCCC},11) \end{array} \right\}
 \end{aligned}$$

$$\text{Weighted Jaccard } J_d^w(x, y) = 1 - \frac{k+2}{n} \quad \text{Edit distance } E_d(x, y) \geq 1 - \frac{2k}{n}$$

Weighted Jaccard = Jaccard for multi-sets

Jaccard and weighted Jaccard ignore relative order

$x = \text{CCCCACCAACACAAAACCC}$

$y = \text{AAAACACAACCCCAACAAA}$

Jaccard and weighted Jaccard ignore relative order

$$\begin{aligned}x &= \text{CCCCACCAACACAAAACCC} && \rightarrow \left\{ \begin{array}{l} \text{AAAA, AAAC, AAC A, AACC, ACAA, ACAC, ACCA, ACCC} \\ \text{CAAA, CAAC, CACA, CACC, CCAA, CCAC, CCCA, CCCC} \end{array} \right\} \\y &= \text{AAAACACAACCCCAACAAA} && \rightarrow \left\{ \begin{array}{l} \text{AAAA, AAAC, AAC A, AACC, ACAA, ACAC, ACCA, ACCC} \\ \text{CAAA, CAAC, CACA, CACC, CCAA, CCAC, CCCA, CCCC} \end{array} \right\}\end{aligned}$$

x, y : de Bruijn sequences,
contain all 16 possible 4-mers once
 $(\sigma!)^{\sigma^{k-1}}$ de Bruijn sequences of length $\sigma^k + \sigma - 1$

Jaccard and weighted Jaccard ignore relative order

$$\begin{aligned}x &= \text{CCCCACCAACACAAAACCC} && \rightarrow \left\{ \begin{array}{l} \text{AAAA, AAAC, AAC A, AAC C, ACAA, ACAC, ACCA, ACCC} \\ \text{CAAA, CAAC, CACA, CACC, CCAA, CCAC, CCCA, CCCC} \end{array} \right\} \\y &= \text{AAAACACAACCCCAACC AAA} && \rightarrow \left\{ \begin{array}{l} \text{AAAA, AAAC, AAC A, AAC C, ACAA, ACAC, ACCA, ACCC} \\ \text{CAAA, CAAC, CACA, CACC, CCAA, CCAC, CCCA, CCCC} \end{array} \right\}\end{aligned}$$

x, y : de Bruijn sequences,
contain all 16 possible 4-mers once
 $(\sigma!)^{\sigma^{k-1}}$ de Bruijn sequences of length $\sigma^k + \sigma - 1$

$$J_d(x, y) = J_d^w(x, y) = 0 \quad E_d(x, y) = 0.63$$

Jaccard is different from edit distance

Unlike edit distance, k -mer Jaccard is insensitive to:

1. k -mer repetitions
 2. relative positions of k -mers
- k -mer Jaccard is not an LSH for the edit distance
 - Still provides big computation saving: asymmetric error model

Jaccard is different from edit distance

Unlike edit distance, k -mer Jaccard is insensitive to:

1. k -mer repetitions
 2. relative positions of k -mers
- k -mer Jaccard is not an LSH for the edit distance
 - Still provides big computation saving: asymmetric error model

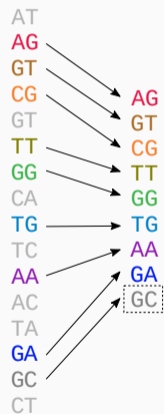
- minHash is an LSH for Jaccard
- OMH is a refinement of minHash
- OMH is sensitive to
 - repeated k -mers
 - relative order of k -mers

$S = \text{AGTTGAGCGGAAGGTG}$, $k = 2$

minHash & OMH sketches

$$S = \text{AGTTGAGCGGAAGGTG}, k = 2$$

π : permutation of Σ^k



minHash & OMH sketches

$S = \text{AGTTGAGCGGAAGGTG}$, $k = 2$, $L = 3$

π : permutation of Σ^k

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

minHash & OMH sketches

$S = \text{AGTTGAGCGGAAGGTG}$, $k = 2$, $L = 3$

π : permutation of Σ^k

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

minHash & OMH sketches

$S = \text{AGTTGAGCGGAAGGTG}$, $k = 2$, $L = 3$

π : permutation of Σ^k

Order: permutation of $\Sigma^k \times \{1, \dots, n\}$

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

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

minHash & OMH sketches

$S = \text{AGTTGAGCGGAAGGTG}$, $k = 2$, $L = 3$, $\ell = 2$

π : permutation of Σ^k

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

1	2	3	4	5	6
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
GT, 13	AG, 5	AG, 5	AA, 10	AG, 5	GT, 13
AA, 10	AG, 11	CG, 7	GT, 13	TT, 2	TT, 2
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

minHash & OMH sketches

$S = \text{AGTTGAGCGGAAGGTG}$, $k = 2$, $L = 3$, $\ell = 2$

π : permutation of Σ^k

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

1	2	3	4	5	6
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
GT, 13	AG, 5	AG, 5	AA, 10	AG, 5	GT, 13
AA, 10	AG, 11	CG, 7	GT, 13	TT, 2	TT, 2
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

minHash & OMH sketches

$S = \text{AGTTGAGCGGAAGGTG}$, $k = 2$, $L = 3$, $\ell = 2$

π : permutation of Σ^k

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

1	2	3	4	5	6
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
GT, 13	AG, 5	AG, 5	AA, 10	AG, 5	GT, 13
AA, 10	AG, 11	CG, 7	GT, 13	TT, 2	TT, 2
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
GC	GA	AG	GG	AG	TG

minHash & OMH sketches

$$S = \text{AGTTGAGCGGAAGGTG}, k = 2, L = 3, \ell = 2$$

Jaccard:

$$\text{Sk}(S) = \begin{pmatrix} \text{GC} \\ \text{TG} \\ \text{GT} \end{pmatrix}$$

OMH:

$$\text{Sk}(S) = \begin{pmatrix} \text{GC} & \text{CA} \\ \text{AG} & \text{GG} \\ \text{AG} & \text{TG} \end{pmatrix}$$

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.

- p_1 : related to probability of hash collisions of weighted Jaccard
- p_2 : related to length of increasing sequence given weighted Jaccard

Jaccard:

- Can use canonical k -mers
- Difficult to find independent hashes:
use bottom sketches ($L \ll n$)

OMH:

- ℓ times as large (cost to encode order)
- $\ell = 1$: LSH / unbiased estimator of weighted Jaccard
- Can't use canonical k -mers: double sketch

OMH has a large gap

- $|S| = 100, k = 5$
- Current proof has a large gap
- What is smallest gap possible?
- OMH/minHash similar to embedding in Hamming space: gap probably unavoidable

