

**University of Stuttgart**

Institute of Biomedical Genetics  
RNA Biology and Bioinformatics

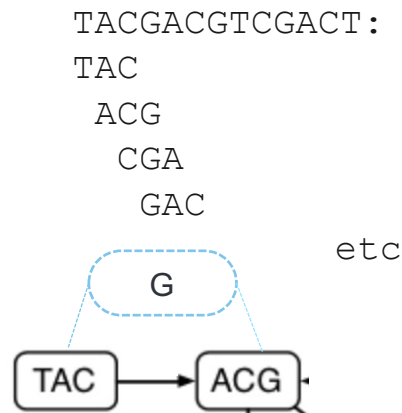
# Succinct de Bruijn graphs

Alex Bowe, Taku Onodera, Kunihiro Sadakane  
and Tetsuyo Shibuya



# De Bruijn Graph

- De Bruijn Graph – DBG
  - cool data structure for DNA sequencing
  - pre-step in DNA assembly
- “De Bruijnizing” algorithm
  1. Chop the reads into **k-mers** (strings of length  $k$ )
  2. Build the **directed graph** based on the overlaps
  3. Call the k-mers **nodes**
  4. Call the overlaps **edges**



# Storage of DBG

- How would you normally store DBG
  - Naive implementation: dictionary
    - Pros: Can be easily controlled/efficient **query operations**
    - Cons: **A lot of space** especially in case of metagenomics

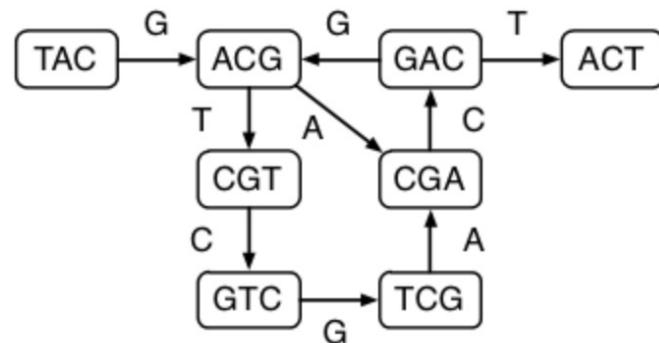
- Optimize space

## ❑ Rules:

- Keep efficient query operations
- Compression allowed while saving the k-mers

## ✓ Solution:

- Information-theoretic lower bound -  $Z$
- Or close to  $Z$  (let me explain)



From	To
TAC	ACG
ACG	CGT, CGA
CGT	GTC
GTC	TCG
TCG	CGA
CGA	GAC
GAC	ACT

# Information Theoretic Lower Bound - OPT

- Given: information-theoretic lower bound:  $Z$
- 3 ways to keep efficient query operations:
  - Implicit:  $Z + \mathcal{O}(1)$  bits
  - **Succinct:  $Z + o(Z)$  bits**
  - Compact:  $\mathcal{O}(Z)$  bits
- *Succinct data structures:*
  - use space “close” to theoretical lower bound
  - allow efficient query operations
  - with efficient “dynamicity” also allow addition/removal of nodes

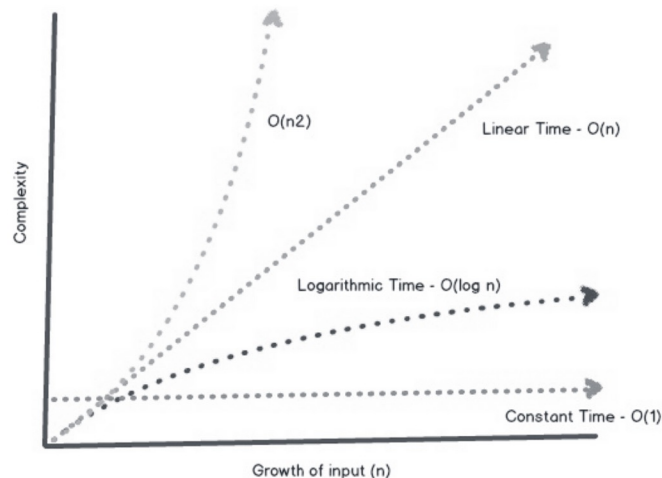
*Information Theoretic Lower Bound:*

Round 6.449 to one character after .

Rule: Max rounding would be 0.099.

$$Z = 6.4$$

$o(Z)$  unit is any function that grows slower than a linear function, e.g.:  $\sqrt{Z}$ ,  $\lg(Z)$ , etc.



# Succinct de Bruijn Graphs

Alexander Bowe<sup>1</sup>, Taku Onodera<sup>2</sup>, Kunihiro Sadakane<sup>1</sup>, and Tetsuo Shibuya<sup>2</sup>

<sup>1</sup> National Institute of Informatics, 2-1-2 Hitotsubashi, Chiyoda-ku,  
Tokyo 101-8430, Japan  
`{alex,sada}@nii.ac.jp`

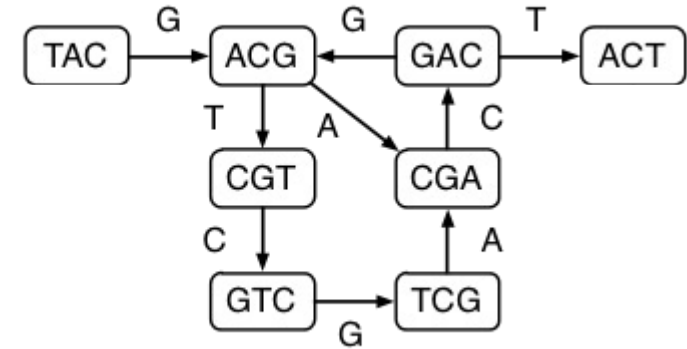
<sup>2</sup> Human Genome Center, Institute of Medical Science,  
University of Tokyo 4-6-1 Shirokanedai, Minato-ku, Tokyo 108-8639, Japan  
`{tk-ono,tshibuya}@hgc.jp`

# SDBG – Succinct de Bruijn Graph

- Goal:
  - represent de Bruijn graph succinctly  $Z + o(Z)$
  - enable efficient query operations
- Succinct de Bruijn graph:
  - DNA sequence of length:  $N$
  - number of edges:  $m$
  - $4m + o(m)$
- Comparisons:
  - distributed hash table(compressed dictionary): 336 GB
  - sparse bitvector representation:
    - $\mathcal{O}(mk)$  w/ encoding and compression: 32 GB
  - ***succinct representation***: 2.5 GB

# Construction

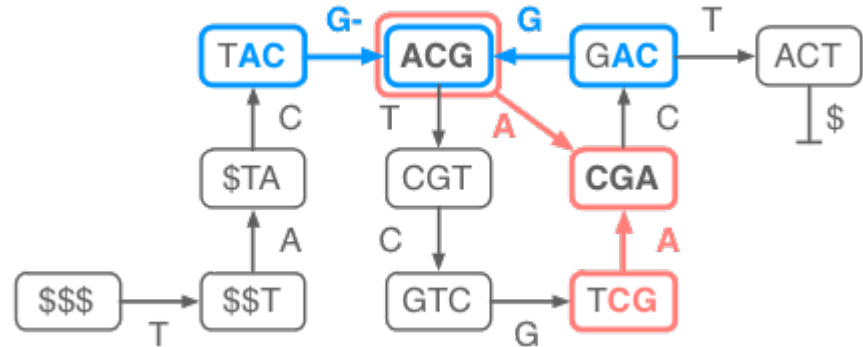
- Take every  $\langle node, edge \rangle$
- Sort them based on reverse(colex order) of node label
- Remove duplicates
- Additional steps:
  - Use padding to know the start and end
  - Rules for padding:
    - $k$  "\$"s to the beginning
    - 1\$ to the end



<i>Node</i>	<i>W</i>
C G A	C
G A C	G
G A C	T
T A C	G-
G T C	G
A C G	A
A C G	T
T C G	A-
C G T	C

# Construction

- $W$  – edge vector(1 char)
- disambiguate identically labelled incoming edges:
  - $A$  and  $A-$ ,  $G$  and  $G-$
- each outgoing edge is stored contiguously
- Bit-vector  $L$ :
  - introduce a bit vector(0 or 1):
  - represent whether an edge is the *last* edge exiting a node
- Bit-vector  $F$ :
  - indicates the first and last occurrences of characters from alphabet  $\{\$, A, C, G, T\}$
  - works due to colex ordering



	$i$	$L$	Node	$W$
	0	1	\$ \$ \$	T
	1	1	\$ C G	C
	2	1	\$ \$ T	C
	3	0	G A C	G
	4	1	G A C	T
	5	1	T T C	G
	6	1	G T A	G
	7	0	A C G	A
	8	1	A C G	T
	9	1	T C G	A-
	10	1	\$ \$ T	A
	11	1	A C T	\$
	12	1	C G T	C



# Rank and Select

*RANK and SELECT have been used before on binary alphabets.*

## RANK

	\$	A	A	C	C	C	C	G	G	G	T	T
F	0	1	2	3	4	5	6	7	8	9	10	11

- Returns the number of elements equal to  $q$  up to position  $x$

$$\text{rank}_q(x) = |\{l, l \in [0 \dots x], S[l] = q\}|$$

$$\text{rank}_A(7) = 2$$

(num of **A**'s are there until position 7)

$$\text{rank}_C(5) = 3$$

(num of **C**'s are there until position 5)

## SELECT

	\$	A	A	C	C	C	C	G	G	G	T	T
F	0	1	2	3	4	5	6	7	8	9	10	11

- Returns the position of the  $x^{\text{th}}$  occurrence of  $q$ .

$$\text{select}_q(x) = \min\{l \in [0 \dots x] : \text{rank}_q(l) = x\}$$

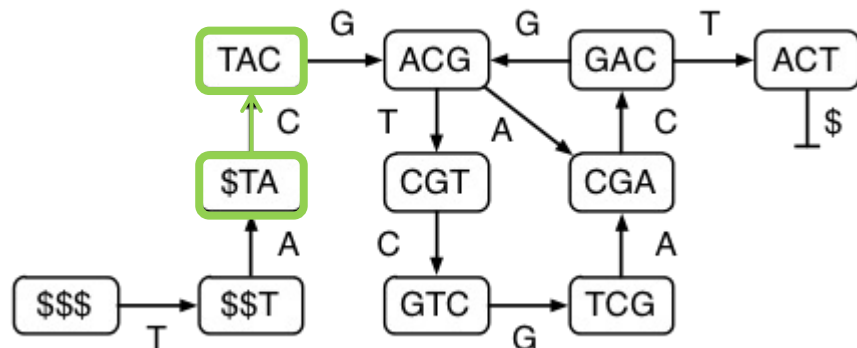
$$\text{select}_T(1) = 10$$

(position of 1'st occurrence of **T**)

$$\text{select}_G(3) = 9$$

(position of 3'rd occurrence of **G**)

$$\text{forward}(v, c) = w$$



	$F$	$i$	$L$	Node	$W$
		0	1	\$ \$ \$	T
		1	1	C C G A	C
		2	1	\$ T A	C
		3	0	G A C	G
		4	1	G A C	T
		5	1	T A C	G
		6	1	G T C	G
		7	0	A C G	A
		8	1	A C G	T
		9	1	T C G	A
		10	1	\$ \$ T	A
		11	1	A C T	\$
		12	1	C G T	C

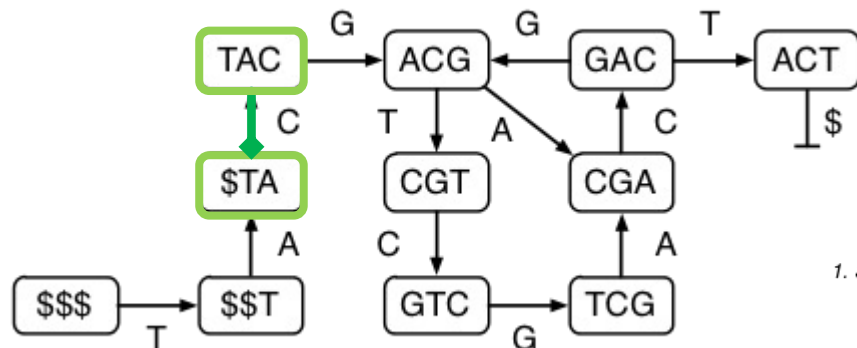
- Find second occurrence of **C** in *Node* column
- Ignore the first **C** since that has  $L$  equal 0
- $O(\text{forward}) = 1$  time unit
  1.  $W$  access -  $O(1)$
  2. rank over  $W$  -  $O(1)$
  3. rank and select over  $L$  -  $O(1)$
  4. accessing  $F$  -  $O(1)$

2. Starting position of  $C = 3$

3. Rank to base = 3

4. Select to  $(3 + 2)$ th position (i.e. the **second C**)

$backward(w, c) = v$  (*inverse\_forward*)



1. Starting position  
of  $C = 3$

$F$	$i$	$L$	Node	$W$
\$ 0	1	3	\$ \$ \$	T
A 1	1	1	C G A	C
C 3	1	1	\$ T A	C
G 7	0	0	G A C	G
T 10	1	1	G A C	T
	5	1	T A C	G
	6	1	G T C	G
	7	0	A C G	A
	8	1	A C G	T
	9	1	T C G	A
	10	1	\$ \$ T	A
	11	1	A C T	\$
	12	1	C G T	C

2. Rank to base = 3

3.1. Rank to current edge = 5,

3.2.  $5 - 3 = 2$ , so we are at the **second C**

4.  $select_C(2) = 2$

- Find second occurrence of **C** in  $W$  column
- We ignore  $L[index(C)] == 0$
- $O(backward) = 1$  time unit
  - Node access -  $O(1)$
  - rank over  $Node$  -  $O(1)$
  - rank and select over  $L$  -  $O(1)$
  - accessing  $W$  -  $O(1)$

$label(v) = dna\_sequence$

$i$	$L$	Node	$W$
0	1	\$	T
1	1	A	C
2	1	A	C
3	0	C	G
4	1	C	T
5	1	C	G
6	1	C	G
7	0	G	A
8	1	G	T
9	1	G	A
10	1	T	A
11	1	T	\$
12	1	T	C

- What if we need to print out the sequence of Node(e.g. *node\_8*)
- Execute backward function  $k$  times
- Function backward takes  $\mathcal{O}(1)$  time
- Repeat the backward function  $k$  times  $\rightarrow$

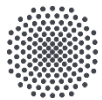
$\mathcal{O}(k)$

# Operations and Time Complexity

Operation	Description	Complexity
$outdegree(v)$	Return number of outgoing edges from node $v$ .	$\mathcal{O}(1)$
$outgoing(v, c)$	From node $v$ , follow the edge labeled by symbol $c$ .	$\mathcal{O}(1)$
$label(v)$	Return (string) label of node $v$ .	$\mathcal{O}(k)$
$indegree(v)$	Return number of incoming edges to node $v$ .	$\mathcal{O}(1)$

## Lesson and Conclusion

- Using “creative counting” with two operations *rank* and *select* the authors created a memory efficient data structure
- The total space has been reduced to the “closer” to OPT space units
  - i.e.: 300 GB → 2.5 GB
- In principle fairly straight forward to implement
- Tools and libraries(C++ based) like: [MegaHIT](#), [BOSS](#), [DynaBOSS](#) and [cosmo](#) have already implemented the data structure
  - For [MegaHIT](#) it took 12 hrs to build a graph from 100 GB Metagenomic reads
  - For slightly further future we might use [DynaBOSS](#)



University of Stuttgart



**Thank you!**