

# Rapid Phylogenies

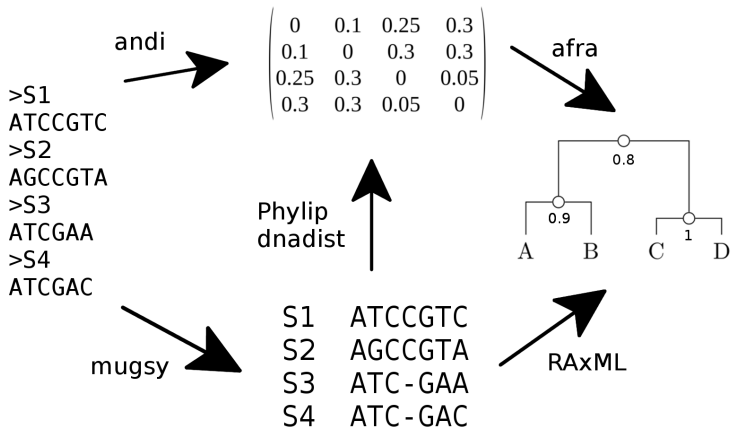
Fabian Klötzl

MPI for Evolutionary Biology, Plön

Hamburg, 2016-05-12

<http://kloetzl.info/downloads/slides-hamburg-2016-05-12.pdf>

# You Had One Job



# Our Plan for Today

- 1 Estimating Evolutionary Distances
- 2 Implementing an Enhanced Suffix Array
- 3 Alignment-Free Support Values
- 4 Fast Multiple Sequence Alignment

## Section 1

# Estimating Evolutionary Distances

# Alignment

$S_1$ : A A G T T A G T G A A C C

$S_2$ : A A G T C A T T G A A A C C

# Alignment

$S_1$ :	A	A	G	T	T	A	G	T	G	A	A	-	C	C
$S_2$ :	A	A	G	T	C	A	T	T	G	A	A	A	C	C

# Alignment

$S_1$ :	A	A	G	T	T	A	G	T	G	A	A	-	C	C
$S_2$ :	A	A	G	T	C	A	T	T	G	A	A	A	C	C

$$d(S_1, S_2) \approx \frac{2}{14}$$

# Alignment

$S_1$ :	A	A	G	T	T	A	G	T	G	A	A	-	C	C
$S_2$ :	A	A	G	T	C	A	T	T	G	A	A	A	C	C

$$d(S_1, S_2) \approx \frac{2}{14}$$

How does one deal with indels?



# Alignment

$S_1$ :	A	A	G	T	T	A	G	T	G	A	A	-	C	C
$S_2$ :	A	A	G	T	C	A	T	T	G	A	A	-	C	C

$$d(S_1, S_2) = \frac{2}{13}$$

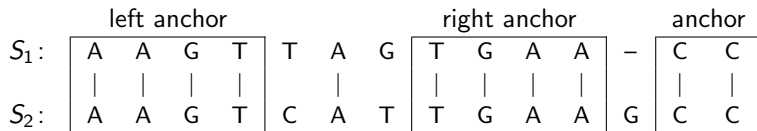
Indels are ignored for the evolutionary distance in common models (Jukes-Cantor, Kimura 2 parameter)!

If indels can be ignored, let's not compute them!

# A Blast from the Past



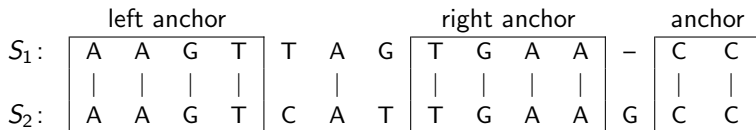
# A Blast from the Past



## Anchor Pair

Two anchors form a pair if they are *equidistant*.

# A Blast from the Past



## Anchor Pair

Two anchors form a pair if they are *equidistant*.

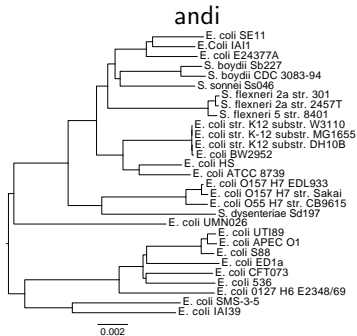
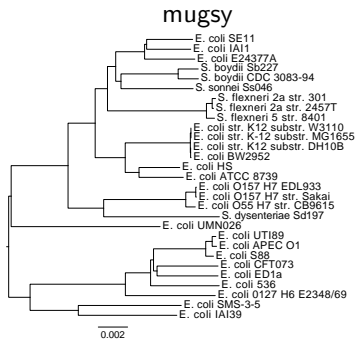
## Anchor Distance

$$d_a = \frac{\#SNP}{\#HomolNucl} = \frac{2}{11}$$

`https://evolbioinf.github.io/andi/`

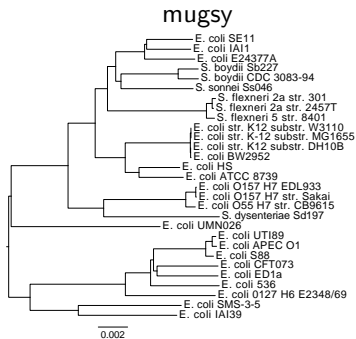
# Demo

# 29 whole *E. Coli*/ *Shigella* Genomes

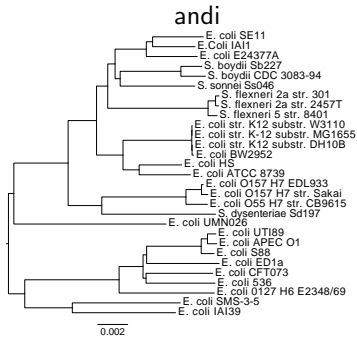




# 29 whole *E. Coli*/ *Shigella* Genomes



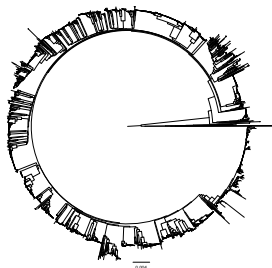
10 164 s



??? s

# Size Matters

- 3085 *Streptococcus pneumoniae* genomes
- multiple contigs per genome
- 2 million nucleotides per genome
- a total filesize of 6.8 GB
- over 9 million pairwise comparisons
- cophylog: 36 days, 2.3 GB
- andi v0.10: 4 h 59 min, 9.8 GB
- most tools exceed 256 GB memory



## Section 2

# Implementing an Enhanced Suffix Array

# getInterval with Child Tables

Not all  $O(1)$  are created equal

```
1  input ( $l-[i..j]$ ,  $m$ ),  $a$ 
2
3  do
4      if  $T[SA[m] + l] = a$  then
5           $j \leftarrow m - 1$ 
6           $m \leftarrow CLD[j+1].L$ 
7          goto line 17
8      end
9
10     if  $m = j$  then
11         break
12     end
13
14      $m \leftarrow CLD[m].R$ 
15 while  $LCP[m]=l$ 
16 if  $T[SA[l] + l] = a$  then
17      $l \leftarrow LCP[m]$ 
18     output ( $l-[i..j]$ ,  $m$ )
19 else
20     output  $\perp$ 
21 end
```



Three step recipe to improve performance:

- 1 Measure!
- 2 Do less work
- 3 Do it faster

CppCon 2014: Chandler Carruth "Efficiency with Algorithms, Performance with Data Structures"

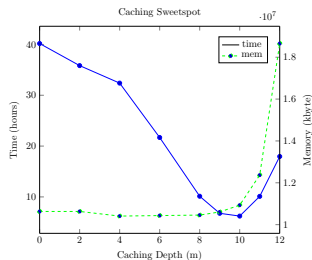
## Measure

# Step 2: Do less work

## Less Memory Accesses

```
amples: 705K of event 'cycles:ppp', Event count (approx.): 410864633684
overhead Command Shared Object Symbol
11,43% andi andi [.] get_interval
10,17% andi andi [.] RMQ_n_1_improved::query
8,89% andi andi [.] esa_init
7,09% andi andi [.] get_match_cached
4,97% andi libdivsufsort.so.3.0.0 [.] divsufsort
4,86% andi andi [.] get_match_from
4,36% andi andi [.] RMQ_nlogn_1::query
1,66% andi andi [.] RMQ_nlogn_1::RMQ_nlogn_1
1,45% andi andi [.] RMQ_n_1_improved::RMQ_n_1_improved
1,11% andi libdivsufsort.so.3.0.0 [.] 0x00000000000061ec
0,99% andi libdivsufsort.so.3.0.0 [.] 0x0000000000006220
0,78% andi andi [.] revcomp
0,67% andi andi [.] dist_anchor
0,53% andi [kernel.vmlinux] [k] pagelock_pfn_to_page
0,49% andi [kernel.vmlinux] [k] clear_page_c_e
0,39% andi libdivsufsort.so.3.0.0 [.] 0x00000000000011a5
0,36% andi [wl] [k] osL_readw
0,36% andi libdivsufsort.so.3.0.0 [.] 0x00000000000011ad
```

- Avoid calling *getInterval*
- Flatten the interval tree into buckets
- Lookup prefix of length  $m$  in  $O(1)$



## Step 3: Do it faster

```
1,28 38: lea    -0x10(%rax),%r15d
        }
        } if( i -- j ){
0,19   cmp    %r15d,%r14d
0,07   + je   %b          break; // 'a' not found, exit early
        }
        // find the next LCP boundary
        m = rmq_lcp->query(i+1, j);
0,39  011: mov    (%r15),%eax
0,14   lea    0x1(%r14),%esi
0,01   mov    %r12,%rdi
0,01   mov    %r15,%edx
1,42   + callq *(%rax)
        } while( LCP[m] == l);
0,48   mov    0x0(%rsp),%rdi
0,40   movsilo %eax,%rdi
2,31   cmp    (%rdi,%r14,4),%r15d
2,33   + jne  %b
        /* We now use an abstract binary search. Think of 'i' as the
        * lower and 'j' as the upper boundary. Then 'm' is the new
        * middle. */
        do {
        // check if 'm' will be a new lower or upper bound.
        if( S[ SA[m] + 1 ] <= a ){
14,07  7f: mov    0x1(%b,%r14,4),%edx
0,36   movzbl 0x1b(%rsp),%ecx
0,21   add    %r13d,%edx
0,32   movsilo %edx,%rdx
1,16   + call *(%eax,%rdx,1),%eax
1,18   + [] %b
0,57   mov    %eax,%r14d
        } else {
        } j = m - 1;
        }
        }
```

- $T[ SA[m] + 1 ]$
- two memory lookups
- unpredictable
- not in cache
- only uses one byte of a whole cache line



## Step 3: Do it faster

```
1,28 38: lea    -0x10(%rax,%r15)
          }
          if( i == j ){
0,19    cmp    %r15,%r14
0,07    + je   %b          break; // 'a' not found, exit early
          }
          // find the next LCP boundary
          m = rmq_lcp->query(i+1, j);
0,39  011: mov    %r12,%rax
0,14    lea    0x1(%r14),%esi
0,01    mov    %r12,%rdi
0,01    mov    %r15,%edx
1,42    + callq *(%rax)
          } while( LCP[m] == l );
0,48    mov    0x0(%rsp),%rdi
0,40    movsllq %eax,%rsi
2,31    cmp    (%rdi,%rsi,4),%r13d
2,33    + jne   %b
          /* We now use an abstract binary search. Think of 'i' as the
          * lower and 'j' as the upper boundary. Then 'm' is the new
          * middle. */
          do {
          // check if 'm' will be a new lower or upper bound.
          if( S[ SA[m] + l ] <= a ){
14,07  7f: mov    0x10(%rbp,%rax,4),%eax
0,36    movzbl 0x1b(%rsp),%ecx
0,21    add   %r13d,%edx
0,32    movsllq %edx,%rdx
1,44    + call *(%rax,%rdx,1),%al
12,18    + {} %a
0,57    mov    %eax,%r14d
          i = m;
          } else {
          j = m - 1;
          }
          }
```

- $\mathcal{T}[ SA[m] + l ]$
- two memory lookups
- unpredictable
- not in cache
- only uses one byte of a whole cache line

Fortunately,

- $l$  is known ahead
- Precompute  $\mathcal{T}[ SA[m] + LCP[m] ]$

# The First Variant Character

Yet another TLA

$i$	$SA$	$LCP$	$T^{SA[i]}$	lcp – intervals		
0	4	-1	AAGG	0	1	3
1	0	3	<u>AAG</u> TAAGG			2
2	5	1	<u>AGG</u>		1	
3	1	2	<u>AG</u> TAAGG			
4	7	0	<b>G</b>			
5	6	1	<u>GG</u>			
6	2	1	<u>G</u> TAAGG			
7	3	0	<b>T</b> AAGG			
8		-1				

# getInterval with Child Tables and FVC

```
1   $ii \leftarrow i$ 
2   $c \leftarrow T[SA[m] + l]$ 
3  goto line 7
4  do
5       $c \leftarrow FVC[j]$ 
6
7      if  $c = a$  then
8           $n \leftarrow CLD[j+1].L$ 
9          output ( $LCP[n] - [i..m-1], n$ )
10     end
11
12      $i \leftarrow m$ 
13     if  $i = j$  then
14         break
15     end
16
17      $m \leftarrow CLD[m].R$ 
18 while  $LCP[m]=l$ 
19 if  $ii = i$  then
20      $c \leftarrow T[SA[j] + l]$ 
21 else
22      $c \leftarrow FVC[j]$ 
23 end
24 if  $c = a$  then
25      $l \leftarrow LCP[m]$ 
26     output ( $l - [i..j], m$ )
27 else
28     output  $\perp$ 
29 end
```

- 20% speed-up
- DC-array by Wu (2016)
- Optimistic algorithm

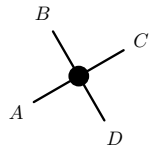
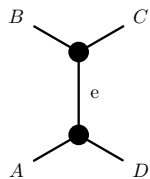
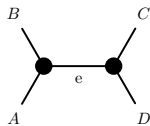
## Section 3

# Alignment-Free Support Values

# Quiz

4

A	0	0.3	0.2	0.1
B	0.3	0	0.1	0.3
C	0.2	0.1	0	0.3
D	0.1	0.3	0.3	0



# Quartet Support Values

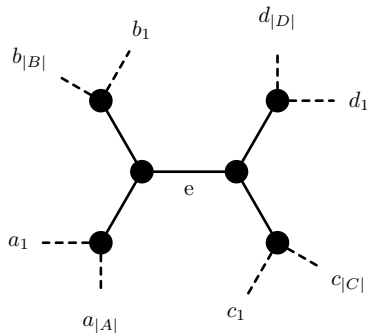
No alignment  $\Rightarrow$  no bootstrapping!

## Quartets

A quartet  $a, b, c, d$  is called supporting an edge  $e$  if

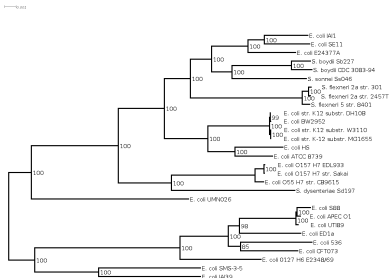
$$D(a, b) + D(c, d) < \min\{D(a, c) + D(b, d), D(b, c) + D(a, d)\}.$$

$$SV_e = \frac{\sum_{a,b,c,d} \mathbf{1}(a,b,c,d \text{ support } e)}{|A||B||C||D|}$$



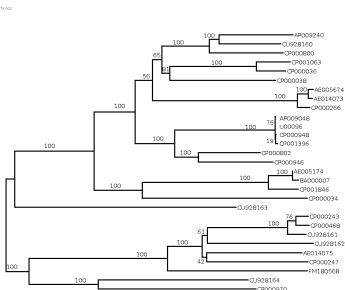
# Eco29 again

mugsy + RAxML



2h 49min + 6h

andi + afra



27s + 1s

## Section 4

# Fast Multiple Sequence Alignment



# Align to Reference



## Demo

# Thank you for your attention



Bernhard Haubold

- `apt-get install andi`
- `brew install science/andi`
- `aura -A andi`
- All projects: [github.com/evolbioinf](https://github.com/evolbioinf)