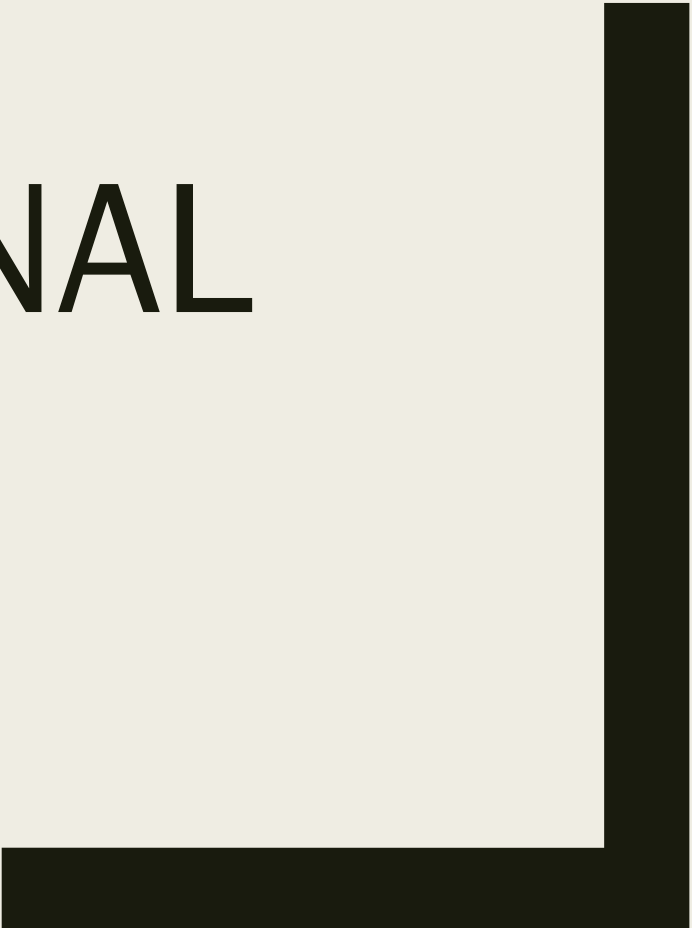# CS 364 COMPUTATIONAL BIOLOGY

Sara Mathieson

Haverford College
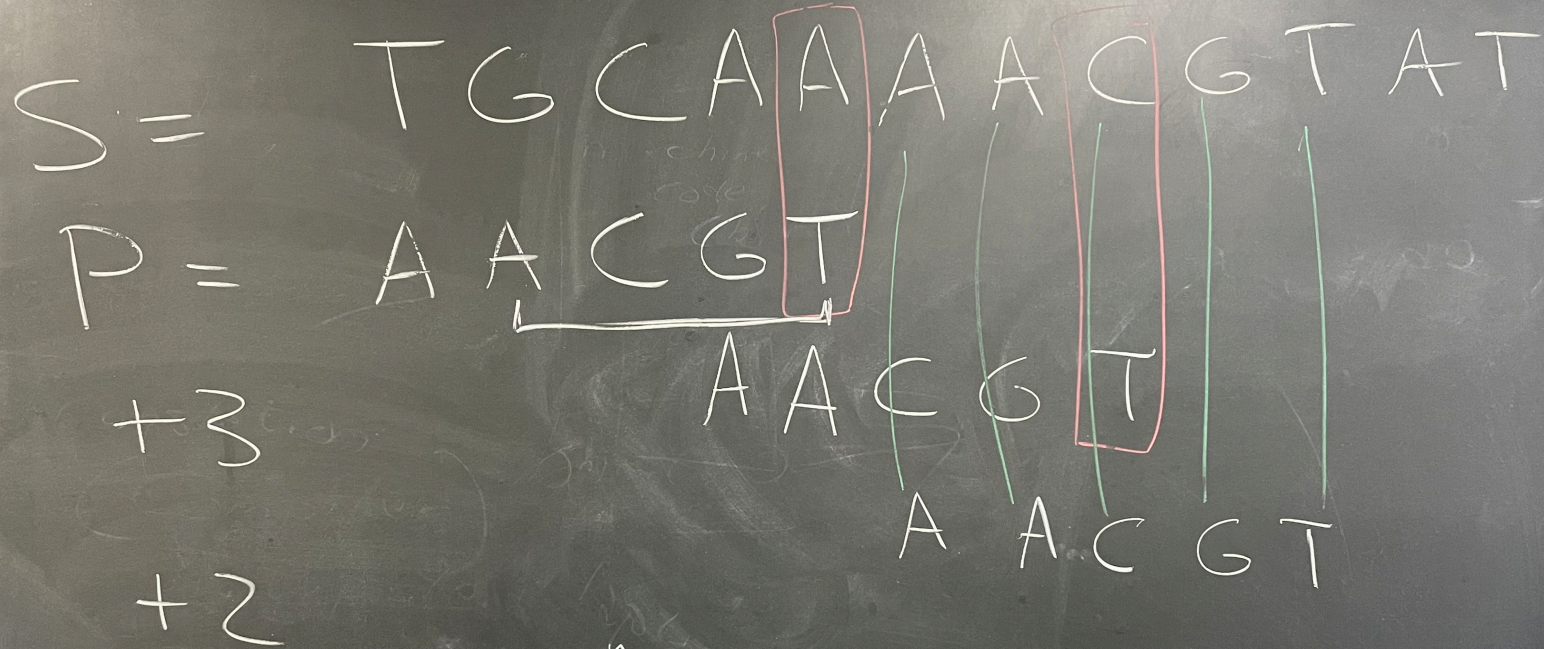
# Outline

- ■ String search: finish Boyer-Moore algorithm

- ■ Sequencing pipeline overview

- ■ Read mapping

- ■ Burrows-Wheeler Transform (BWT)

# Boyer-Moore Algorithm

Bad char worksheet

S = TGCAAAAACGTAT

P = AACGT

+3    AACGT

+2    AACGT

Num compare = 3 + 5 = 8
                ↑       ↑
            mismatch  match

```
S = AAATAAATAAATAAAT
    AAAA
        AAAA
```

`skip`

```
S = AAAAAAAAAAAATAAA
    TAAA
      TAAA
```

`No skip`

```
S = AAAAAAAAAAATAAA
    AAAT
       AAAT
```

`No skip`

Good Case

Bad Cases

*Worst-case* time complexity proportional to $nm$, where $n$ is size of S and $m$ is size of the pattern, P.   However, in practice *average* time complexity is very good.

Can we do better in these worst case situations?

Idea: there is information in "matches so far"

S = AAAAAAAAAAAATAAA
    TAAA

Found a mismatch at the left most position, where the pattern has a T but we know that S has AAA to the right of that, so we could never match in that region and can skip over it.

S = AAAAAAAAAAAAAAA
+4 ⟶ TAAA

**Good Suffix Rule:** Consider all the possible points in P where mismatch can happen. For each kind of mismatch, look left and see if the mismatch pattern up to the current point can be found in P again.

```
  123456789
S:TAT**T**CGGTT        **^ACG**
P:GCG**A**CG
```

```
  123456789
S:TAT**T**CGGTT
P:    GCG**A**CG
```

Notation: We use ^X to mean "not X"

So ^A = C, G or T

^ACG = CCG, GCG or TCG etc…

```
        123456789
S:TATTCGGTT
P:GCGACG
```

```
  123456789
S:##^ACG###
P:GCGACG
```

```
   123456789
S:##^ACG###
P: GCGACG        +1
```

```
     123456789
S:##^ACG###
P:  GCGACG        +2
```

```
   123456789
S:##^ACG###
P:   GCGACG       +3
```

```
   123456789
S:##^ACG###
P:   GCGACG        +3
```

# All possible mismatch patterns for P = "GCGACG"

| Position of first mismatch of P, scanning from right to left | Implied pattern in S | Action |
|---|---|---|
| 6 | ^G | move 1 |
| 5 | ^CG | move 5 |
| 4 | ^ACG | move 3 |
| 3 | ^GACG | move 5 |
| 2 | ^CGACG | move 5 |
| 1 | ^GCGACG | move 5 |

$$\overset{0}{G}\overset{1}{C}\overset{2}{G}\overset{3}{A}\overset{4}{C}\overset{5}{G}$$

| pattern in S | i | shift |
|---|---|---|
| ^G | 5 | 1 |
| ^CG | 4 | 5 |
| ^ACG | 3 | 3 |
| ^GACG | 2 | 5 |
| ^CGACG | 1 | 5 |
| ^GCGACG | 0 | 5 |

good suffix table

# Exercise: P = "TAAAA"

| Position of first mismatch of P, scanning from right to left | Implied pattern in S | Action to be taken |
|---|---|---|
| P[5] | ^A | increase offset by 4 |
| P[4] | ^AA | increase offset by 3 |
| P[3] | ^AAA | increase offset by 2 |
| P[2] | ^AAAA | increase offset by 1 |
| P[1] | ^TAAAA | increase offset by 5 |

# Combine the two rules to complete the Boyer-Moore algorithm

Match P to S, from right to left starting at some offset position of S. If a mismatch is found compute

*shift* = max(bad character rule, good suffix rule)

Increase offset by *shift* and start the matching again from the right most position of P.

# Good suffix rule worksheet

| pattern in S | shift |
|---|---|
| ^T | 1 |
| ^AT | 5 |
| ^GAT | 5 |
| ^GGAT | 5 |
| ^AGGAT | 5 |
| ^TAGGAT | 5 |

good suffix

GAGATTATAGGATTACA
TAGGAT
                TAGGAT        +5
                 TAGGAT       +1
bad char              TAGGAT  +1

GAGATTATAGGATTACA
TAGGAT
            TAGGAT            +4
             TAGGAT          +2
              TAGGAT         +1

cycl
perm

# Boyer-More Complexity

# Best case?

```
S = AAACAAACAAACAAAT
    AAAT
        AAAT
```

Maximum num comparisons: O(N/M)

Gets more efficient as the pattern gets longer!

# Average case?

```
S = AGTCTAGCTAGCATCGACTACGAC
      ACGT
          ACGT
              ACGT
```

Average num comparisons: O(N)

# Small alphabet

```
S = AGTCTAGCTAGCATCGACTACGAC
    ACGT
         ACGT
              ACGT
```

# Large alphabet

```
S = ASDFAJNSDAWZA#XXKXLS#%K
    XZAK
               XZAK
                      AZAK
```

Bigger skip with larger alphabet

# Worst case?

```
S = AAAAAAAAAAAAAAAAAA
    AAAT
       AAAT
```

No skip: O(N)

# Worst case?

```
S:  AAAAAAAAAAAAAAAAAAAAAAA
P:  AAAAA
```

No skip, no mismatch: O(NM)

* Worst case is O(N) if the pattern does not appear in the text

# Remarks on Boyer-Moore

- O(N/M) -> Becomes more efficient, the longer the pattern

- Small memory complexity.

- Not as efficient when we have a small alphabet size

# Remarks on Boyer-Moore

- Boyer-Moore and variations are good general search algorithms

- At the relatively modest cost of O(M) preprocessing, they produce substantial speed-ups over the naïve algorithm

- This is hard to see abstractly, but in practice they have very good performance (i.e. the average case is close to the best case)

- When you hit Ctrl-F in a website or text document, it's running a variation of Boyer-Moore under the hood

# Extra reading

## Fast String Searching

ANDREW HUME
*AT&T Bell Laboratories, 600 Mountain Ave., Murray Hill, NJ 07974, U.S.A.*

AND

DANIEL SUNDAY
*Johns Hopkins University / Applied Physics Laboratory, Johns Hopkins Rd., Laurel, MD 20723, U.S.A.*

### SUMMARY

Since the Boyer-Moore algorithm was described in 1977, it has been the standard benchmark for the practical string search literature. Yet this yardstick compares badly with current practice. We describe two algorithms that perform 47% fewer comparisons and are about 4.5 times faster across a wide range of architectures and compilers.

These new variants are members of a family of algorithms based on the skip loop structure of the preferred, but often neglected, fast form of Boyer-Moore. We present a taxonomy for this family, and describe a toolkit of components that can be used to design an algorithm most appropriate for a given set of requirements.
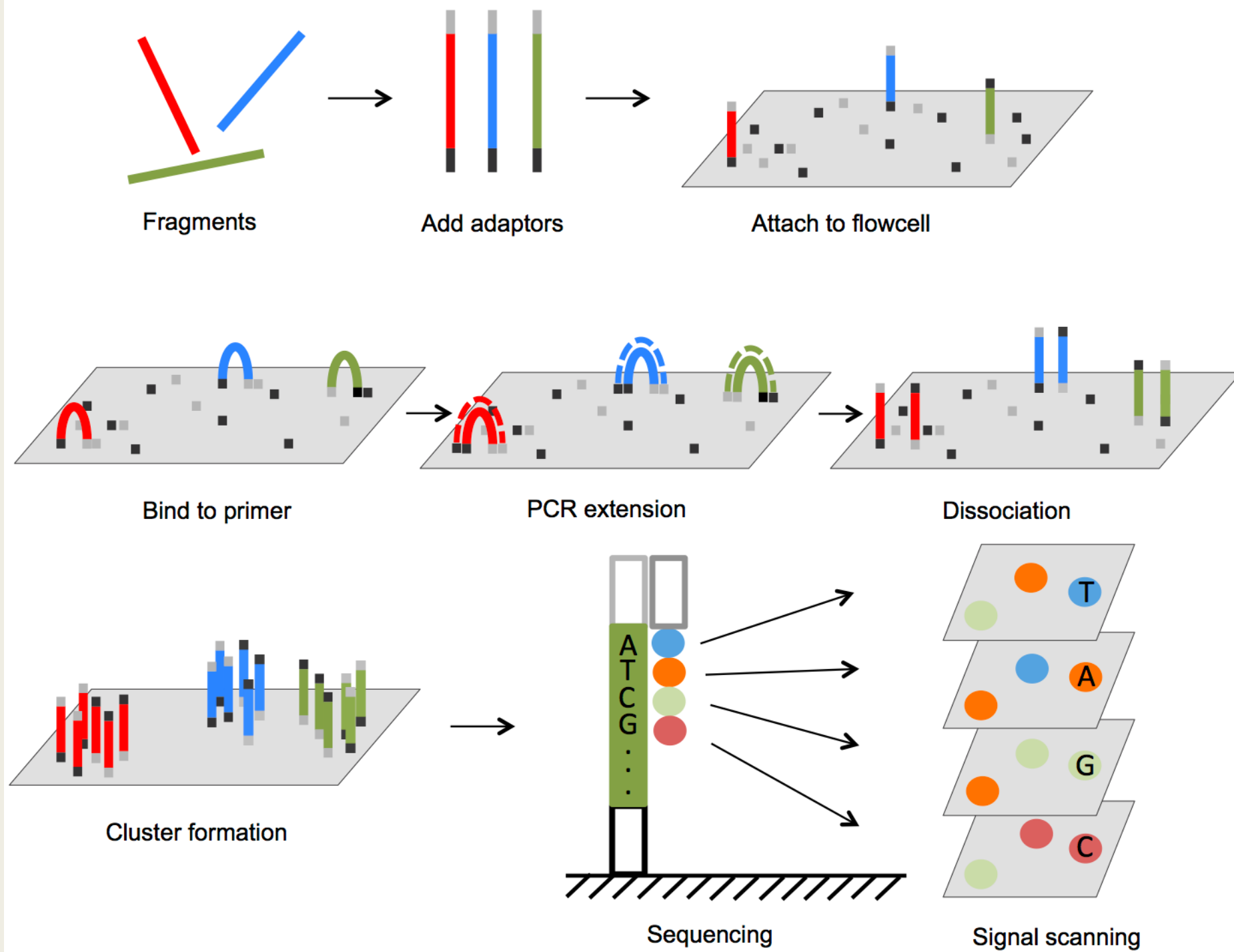
KEY WORDS  String  searching  Pattern  matching  Boyer-Moore

For optimal performance we start to care about 1) constants 2) characteristics of the data

# Final thought…

- We sped up our algorithm by doing some pre-processing of the <u>pattern</u>

- Often (e.g. read mapping), we want to match a large number of patterns (reads) to a search string (reference genome) that doesn't change

- Is there some way we could speed this up by pre-processing the <u>search string</u>?

# Sequencing Pipeline

Fragments · Add adaptors · Attach to flowcell · Bind to primer · PCR extension · Dissociation · Cluster formation · Sequencing · Signal scanning

# Read Mapping

# Read mapping

Long reference genome

…**AGTATCTGTCTTTGATTCCTGCCTCATCCTATTATTTATCGCACCTACGTTCAATATT**…

CTTTGATTCCTGCC                TTATTTATCGCACCTAC

CTGCCTCATCCTA

Many short reads to align

Reference doesn't change so
lets process it once (slow) but
then hopefully it will be fast
to map each read.

# Bowtie and BWA (posted reading)

- First practical read aligners to use the Burrows-Wheeler Transform

**Fast and accurate short read alignment with Burrows−Wheeler transform** 🔓

Heng Li, Richard Durbin ✉     Author Notes

*Bioinformatics*, Volume 25, Issue 14, 15 July 2009, Pages 1754–1760,
https://doi.org/10.1093/bioinformatics/btp324
**Published:** 18 May 2009     **Article history** ▾

BWA

Bowtie

**Ultrafast and memory-efficient alignment of short DNA sequences to the human genome**

Ben Langmead ✉, Cole Trapnell, Mihai Pop and Steven L Salzberg

*Genome Biology* 2009 10:R25
https://doi.org/10.1186/gb-2009-10-3-r25    © Langmead et al.; licensee BioMed Central Ltd. 2009
Received: 21 October 2008 │ Accepted: 4 March 2009 │ Published: 4 March 2009
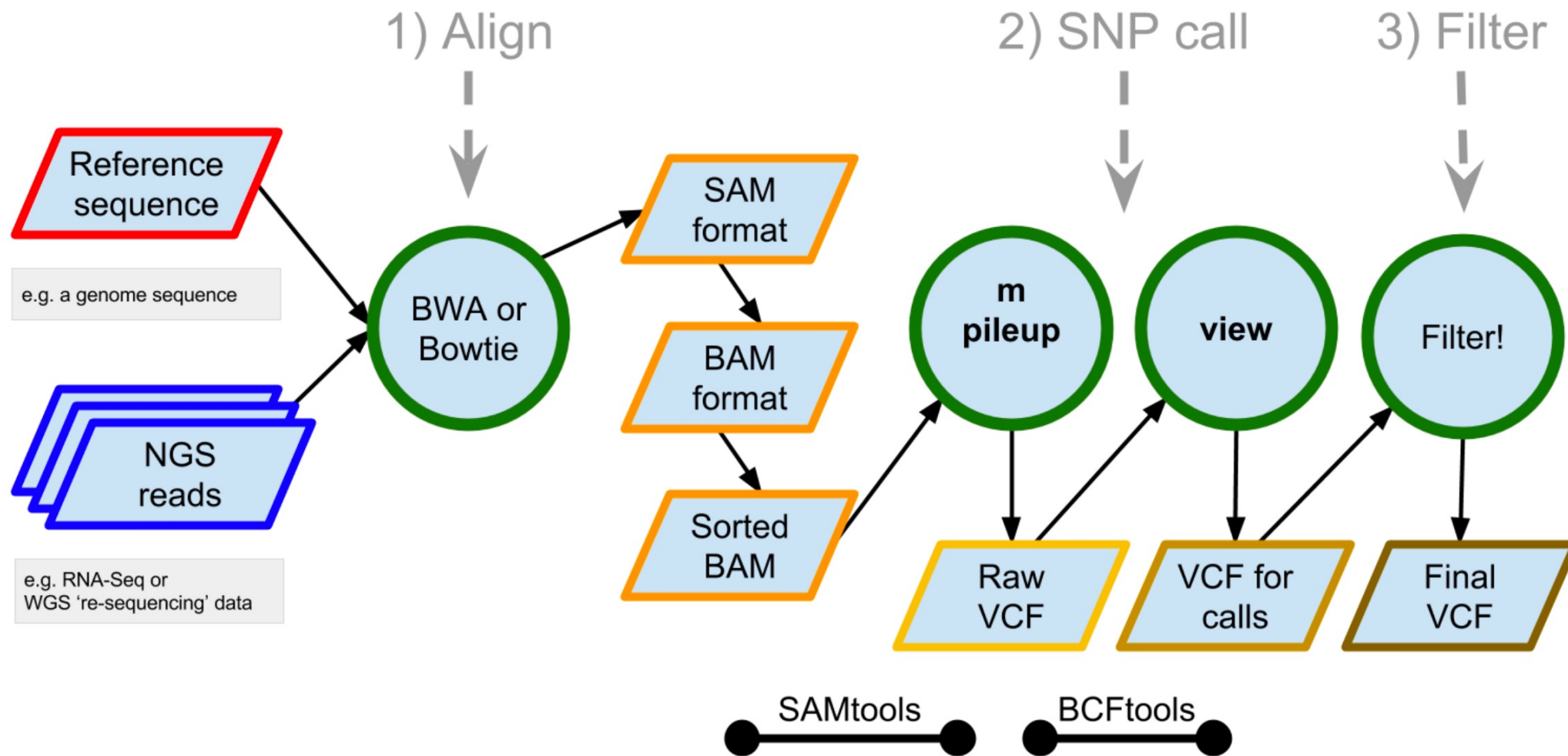
# Comparison of Bowtie and BWA
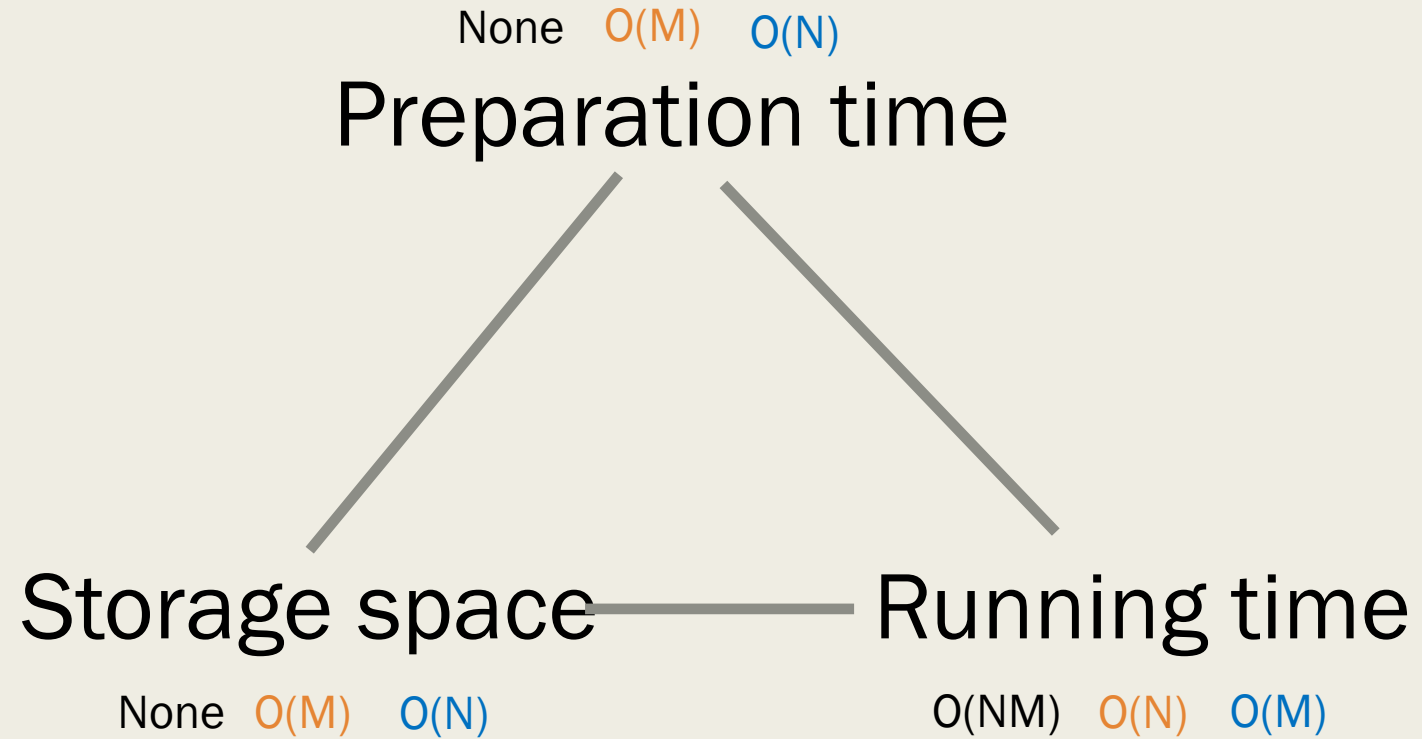
**Table 2.**

Evaluation on real data

| Program | Time (h) | Conf (%) | Paired (%) |
|---------|----------|----------|------------|
| Bowtie | 5.2 | 84.4 | 96.3 |
| BWA | 4.0 | 88.9 | 98.8 |
| MAQ | 94.9 | 86.1 | 98.7 |
| SOAP2 | 3.4 | 88.3 | 97.5 |

The 12.2 million read pairs were mapped to the human genome. CPU time in hours on a single core of a 2.5 GHz Xeon E5420 processor (Time), percent confidently mapped reads (Conf) and percent confident mappings with the mates mapped in the correct orientation and within 300 bp (Paired), are shown in the table.

# Pipeline overview



Credit: Dan Bolster (EBI)

# No such thing as a free lunch

None  O(M)  O(N)
## Preparation time

## Storage space ——— Running time

None  O(M)  O(N)          O(NM)  O(N)  O(M)

Naïve search
Boyer-Moore and variations
Hash tables and BWT

# k-mer hashing

AGTATCTGTCTTTGATTCCTGCCTCATCCTATTATTTATCGCACCTACGTTCAATATT

AGTA: 1

Recall: k-mer is a string of length k. We'll write things like 4-mer, 32-mer etc... to refer to specific lengths

# k-mer hashing

AGTATCTGTCTTTGATTCCTGCCTCATCCTATTATTTATCGCACCTACGTTCAATATT

```
AGTA:  1
GTAT:  2
```

# k-mer hashing

AG TATC TGTCTTTGATTCCTGCCTCATCCTATTATTTATCGCACCTACGTTCAATATT

```
AGTA:  1
GTAT:  2
TATC:  3
```

# k-mer hashing

AGTATCTGTCTGTGATTCCTGCCTCATCCTATTATTTATCGCACCTCTGTTCAATATT

```
AGTA:  1
GTAT:  2
TATC:  3
ATCT:  4
TCTG:  5,9,42
CTGT:  6
TGTC:  7
GTCT:  8
TCTT:  9
  .
  .
  .
  .
  .
```

Collision!

# k-mer hashing

AGTATCT|GTCT|GTGATTCC|TGCCTCATCCTATTATTTATCGCACCTCTGTTCAATATT

AGTA:  1
GTAT:  2
TATC:  3
ATCT:  4
TCTG:  5,9,42
CTGT:  6
TGTC:  7
GTCT:  8
TCTT:  9
.
.
.
.
.

Now, suppose I want to look up:

GTCT|GTGATTCC

1. Take first k-mer in pattern
2. Look up positions in index
3. Check each position for match

How many different k-mers are there?

# Burrows-Wheeler Transform (BWT)

# BWT

$S = \text{banana}\$ \nwarrow$

special char "before" all char

| A |  |  | $\Pi(S)$ | rank | $\Pi^{\text{sorted}}(S)$ ⓕ | ⓛ | BWT | F |
|---|---|---|---|---|---|---|---|---|
| Cycle permutations |  |  | banana$ | 5 | $ banana | | a | ⓢ |
|  |  |  | anana$b | 4 | a $banan | | n } repeat | a |
|  |  |  | nana$ba | 7 | ana$ban | | | a |
|  |  |  | ana$ban | 3 | anana$b | | n | |
|  |  |  | na$bana | 6 | banana$ | | b | b |
|  |  |  | a$banan | 2 | na$bana | | $ | |
|  |  |  | $banana | 1 | nana$ba | | a } repeat | n |
|  |  |  |  |  |  | | a | |

chars

| F | L | backtrace | S | | F | L |
|---|---|---|---|---|---|---|
| $\circled{\$}$ → $a_1$ | | $\$ → a_1$ | $a_1\$$ | | $\$$ | → $e_1$ |
| $a_1$ → $n_1$ | | $a_1 → n_1$ | $na\$$ | | $e_1$ | → $v$ |
| $a_2$ → $n_2$ | | $n_1 → a_2$ | $ana\$$ | | $e_2$ | $r$ $t_1$ |
| $a_3$ → $b_1$ | | $a_2 → n_2$ | $nana\$$ | | $e_3$ | |
| $b_1$ → $\$$ | | $n_2 → a_3$ | $anana\$$ | | $i$ $p$ | |
| $n_1$ → $a_2$ | | $a_3 → b_1$ | $banana\$$ | | $r$ $t_1$ | |
| $n_2$ → $a_3$ | | $b_1 → \$$ | stop | | $t_2$ | |
| | | | | | $v$ | → $e_2$ |