



# TOTALRECALLER: Improved Accuracy and Performance via Integrated Alignment & Base-Calling

Fabian Menges

New York, 2011/3/10



## Overview

### Introduction

### Basecalling with alignment

- Sequencing technology

- Linear error model

- Sequence Alignment

- Branch and Bound

### Results



## Overview

### Introduction

Basecalling with alignment

Sequencing technology

Linear error model

Sequence Alignment

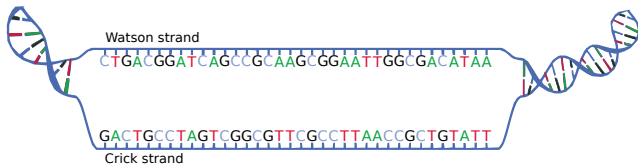
Branch and Bound

### Results



# Introduction

## DNA





## DNA

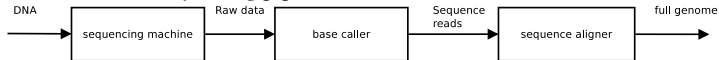
DNA of different organisms:

Organism type	Organism	Genome size (base pairs)
Virus	Bacteriophage MS2	3,569
Virus	SV40	5,224
Bacterium	Haemophilus influenzae	1,830,000
Plant	Populus trichocarpa	480,000,000
Yeast	Saccharomyces cerevisiae	12,100,000
Fungus	Aspergillus nidulans	30,000,000
Insect	Apis mellifera (honey bee)	1,770,000,000
Fish	Tetraodon nigroviridis	385,000,000
Fish	Protopterus aethiopicus	130,000,000,000
Mammal	Homo sapiens	3,200,000,000

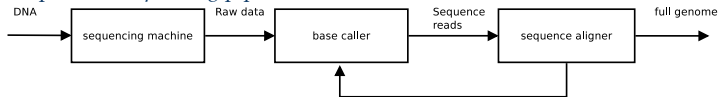


## Re-sequencing Pipeline

### Conventional *re-sequencing* pipeline:



### Proposed *re-sequencing* pipeline:



**Motivation:** Avoiding errors in *sequence reads*.



# Basecalling with alignment

---

## Overview

### Introduction

### **Basecalling with alignment**

Sequencing technology

Linear error model

Sequence Alignment

Branch and Bound

### Results



# Basecalling with alignment

---

## Overview

### Introduction

### Basecalling with alignment

- Sequencing technology

- Linear error model

- Sequence Alignment

- Branch and Bound

### Results





# Basecalling with alignment

## Illumina sequencing machine

*sequence reads* with a length of up to 125BP.

Flow Cell:



Genome Analyzer IIe:



Input: Cluster of DNA-Fragments  
Output: Raw sequence intensities

Source: <http://www.illumina.com>



# Basecalling with alignment

## Illumina Intensities

### Sequencing machine:

- Input: DNA-fragments
- Output: Intensities

Cycle	Channel A	Channel C	Channel G	Channel T
1	15.7	-19.5	<b>3812.9</b>	1398.9
2	-29.0	41.6	365.5	<b>1200.5</b>
3	14.4	36.6	379.1	<b>1447.0</b>
⋮	⋮	⋮	⋮	⋮
76	837.4	549.4	1098.8	841.7
77	633.0	491.8	1280.7	901.4
78	602.9	558.2	1036.2	860.9

### Basecaller:

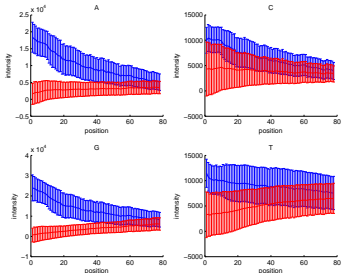
⇒ *sequence read*: GTT...???



# Basecalling with alignment

## Illumina Intensities

Level distance: High-level(blue) , Low-level(red)



⇒ The level distance between low and high signal levels is too small in order to reconstruct the correct sequence read.



# Basecalling with alignment

---

## Overview

### Introduction

### **Basecalling with alignment**

Sequencing technology

**Linear error model**

Sequence Alignment

Branch and Bound

### Results



# Basecalling with alignment

## Causes of Errors

- **Fading:** Signal strength decreases with the number cycles.
- **Crosstalk:** The channels are not independent from one another.
- **Lagging:** There is interference between two consecutive cycles.



# Basecalling with alignment

## Linear error model Crosstalk and fading

- Cycle:  $k \in \mathbb{N}$ :
- **Input:** Raw intensities:  $\mathbf{I}_k = (\mathbf{I}_A^k \quad \mathbf{I}_C^k \quad \mathbf{I}_G^k \quad \mathbf{I}_T^k)^\top$
- **Output:** Filtered intensities:  $\mathbf{X}_k = (\mathbf{X}_A^k \quad \mathbf{X}_C^k \quad \mathbf{X}_G^k \quad \mathbf{X}_T^k)^\top$
- Crosstalk matrix:  $\mathbf{A}_k \in \mathbb{R}^{4 \times 4}$
- Model:  $\mathbf{I}_k = \mathbf{A}_k \cdot \mathbf{X}_k$

$$\mathbf{A}_k = \begin{pmatrix} \mu_{k,A}^A & \mu_{k,C}^A & \mu_{k,G}^A & \mu_{k,T}^A \\ \mu_{k,A}^C & \mu_{k,C}^C & \mu_{k,G}^C & \mu_{k,T}^C \\ \mu_{k,A}^G & \mu_{k,C}^G & \mu_{k,G}^G & \mu_{k,T}^G \\ \mu_{k,A}^T & \mu_{k,C}^T & \mu_{k,G}^T & \mu_{k,T}^T \end{pmatrix}$$

- **Filter:** Filtered intensities:  $\mathbf{X}_k = \mathbf{A}_k^{-1} \cdot \mathbf{I}_k$



# Basecalling with alignment

## Linear error model Crosstalk, fading and lagging

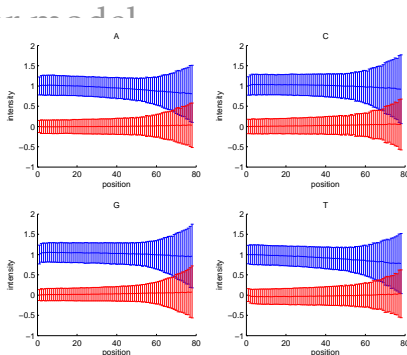
- Cycle:  $k \in \mathbb{N}$ :
- **Input:** Raw intensities:  $\mathbf{I}_k = (\mathbf{I}_A^k \quad \mathbf{I}_C^k \quad \mathbf{I}_G^k \quad \mathbf{I}_T^k)^\top$
- Crosstalk matrix:  $\mathbf{A}_k \in \mathbb{R}^{4 \times 4}$
- Lagging matrix:  $\mathbf{\Upsilon}_k \in \mathbb{R}^{4 \times 4}$
- **Output:** Filtered intensities:

$$\begin{pmatrix} \mathbf{X}_{k-1} \\ \mathbf{X}_k \end{pmatrix} = \underbrace{\begin{pmatrix} \mathbf{A}_{k-1} & 0 \\ \mathbf{\Upsilon}_k & \mathbf{A}_k \end{pmatrix}^{-1}}_{\mathbf{G}_k \in \mathbb{R}^{8 \times 8}} \cdot \begin{pmatrix} \mathbf{I}_{k-1} \\ \mathbf{I}_k \end{pmatrix}$$



# Basecalling with alignment

Linear error  
Result



⇒ The level distance between the high and low level allows to separate the channels up to the 60th cycle.

**Goal:** Correct reconstruction of all bases in a sequence read!





# Basecalling with alignment

---

## Overview

### Introduction

### Basecalling with alignment

Sequencing technology

Linear error model

**Sequence Alignment**

Branch and Bound

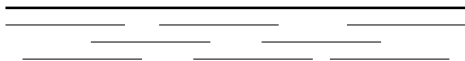
### Results



# Basecalling with alignment

## Sequence Alignment

Reference



Reads

- Discover the correct position for each *sequence reads* in a given reference genome.
- Can be described a search of *sequence reads* in the reference genome considering (*mismatches*).

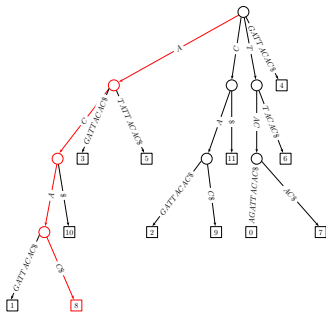
**Problem:** Fast search with mismatches in large reference genome ( $\sim 3 \cdot 10^9$ BP).



# Basecalling with alignment

## Suffix Tree

Example: Reference T = "TACAGATTACAC\$"



- Advantage: Allow to search for a *read* with length  $n$  in  $O(n)$ .
- Disadvantage: Requires for a reference with length  $m$   $O(m^2)$  memory.



# Basecalling with alignment

## Burrows-Wheeler Transformation

- Reversible transformation.
- Was developed for a compression algorithm

Example:

TACAGATTACAC\$		\$TACAGATTACAC	
ACAGATTACAC\$T		AC\$TACAGATTAC	
CAGATTACAC\$TA		ACAC\$TACAGATT	
AGATTACAC\$TAC		ACAGATTACAC\$T	
GATTACAC\$TACA		AGATTACAC\$TAC	
ATTACAC\$TACAG		ATTACAC\$TACAG	
TTACAC\$TACAGA	⇒	C\$TACAGATTACA	⇒
TACAC\$TACAGAT		CAC\$TACAGATTA	CCTTCGAAAAT\$A
ACAC\$TACAGATT		CAGATTACAC\$TA	
CAC\$TACAGATTA		GATTACAC\$TACA	
AC\$TACAGATTAC		TACAC\$TACAGAT	
C\$TACAGATTACA		TACAGATTACAC\$	
\$TACAGATTACAC		TTACAC\$TACAGA	



# Basecalling with alignment

## Ferragina-Manzini search

Allows to search the suffix tree through the BWT.

---

```
Algorithm 1: FM_search
/* Ferragina-Manzini search algorithm */
P search pattern
input : n length of P
        F table of accumulated character frequencies
output: ep - sp + 1 number of occurrences of P in T
i = n - 1;
/* i is a index to the last position of the pattern */
c = P[i]; /* set c to the last character of the pattern */
sp = F[c - 1] + 1;
/* set the start index */
/* c-1 represents the lexicographical previous character */
ep = F[c];
/* set the end index */
while (sp ≤ ep) and (i ≥ 1) do
    i = i - 1;
    /* decrement the position */
    c = P[i];
    /* get the next character */
    sp = F[c - 1] + C(c, sp - 1) + 1;
    /* update the start index */
    ep = F[c - 1] + C(c, ep);
    /* update the end index */
if (ep < sp) then
    return "pattern not found"
else
    return "found (ep - sp + 1) occurrences"
```

---





# Basecalling with alignment

## Overview

### Introduction

### Basecalling with alignment

Sequencing technology

Linear error model

Sequence Alignment

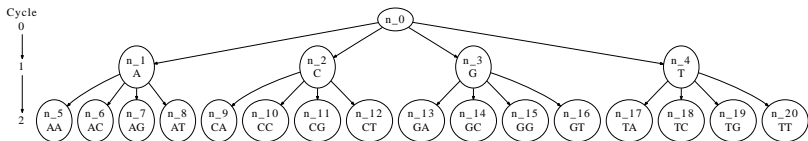
**Branch and Bound**

### Results



# Basecalling with alignment

## Basecalling by building a tree



- Every cycle the tree grows in depth. At cycle  $k$  the tree has  $4^k$  leaves and  $\sum_{i=1}^k 4^i$  nodes.
- Every node represents a possible correct *sequence read*.
- Every node has specific probability of being the correct *sequence read*.
- The node with the highest probability represents the best solution.

**Problem:** Too many nodes! How is the probability for correctness computed?





# Basecalling with alignment

## Branch and Bound

The branch and bound algorithm allows to reduce the number of nodes that need to be considered. Every cycle the following three steps are performed:

- **Branch:** All nodes in the set of possible solutions are expanded.
- **Bound:** The new child nodes are weighted according to their probability for correctness: a score function.
- **Pruning:** Only the best  $b$  nodes are kept in the set of possible solutions.

⇒ Branch and Bound reduces the number of nodes that need to be considered by building the tree only partially.

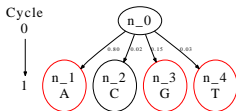
⇒ The correct/best solution can not be guaranteed!



# Basecalling with alignment

## Example

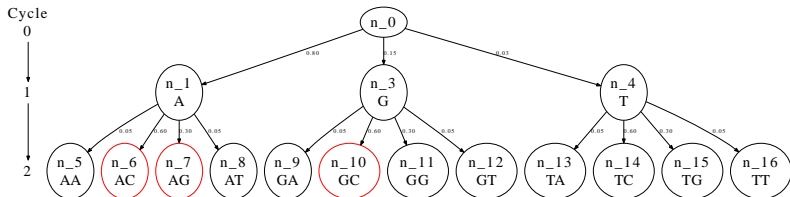
Branch and Bound:  $b = 3$  possible solutions are considered.





# Basecalling with alignment

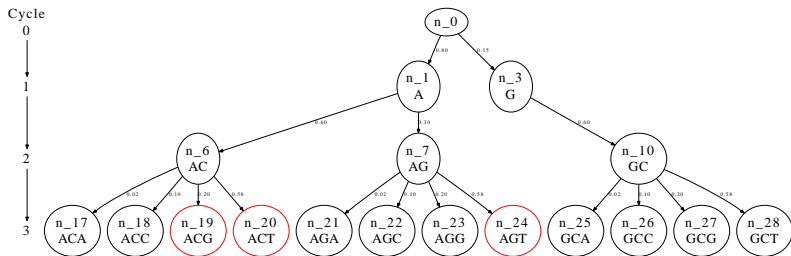
## Example





# Basecalling with alignment

## Example





# Basecalling with alignment

## Score function

$$\begin{aligned} P_k(B | \mathbf{X}_k) &= \frac{P_k(\mathbf{X}_k | B)P_k(B)}{P_k(\mathbf{X}_k)} && \text{with } B \in \{A, C, G, T\} \\ &= \frac{P_k(\mathbf{X}_k | B)P_k(B)}{P_k(\mathbf{X}_k | B)P_k(B) + P_k(\mathbf{X}_k | \neg B)P_k(\neg B)} \\ &= \frac{1}{1 + \frac{P_k(\mathbf{X}_k | \neg B)P_k(\neg B)}{P_k(\mathbf{X}_k | B)P_k(B)}} \\ &= \frac{1}{1 + \underbrace{\frac{P_k(\mathbf{X}_k | \neg B)}{P_k(\mathbf{X}_k | B)}}_{\text{Intensities}} \cdot \underbrace{\frac{P_k(\neg B)}{P_k(B)}}_{\text{Sequence alignment}}} \end{aligned}$$



## Overview

### Introduction

### Basecalling with alignment

- Sequencing technology

- Linear error model

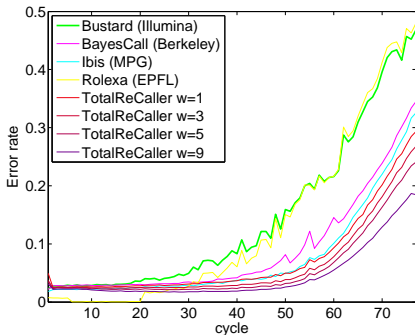
- Sequence Alignment

- Branch and Bound

## Results

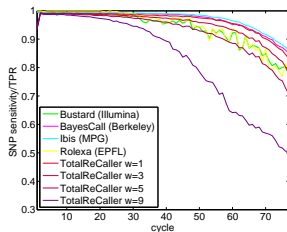
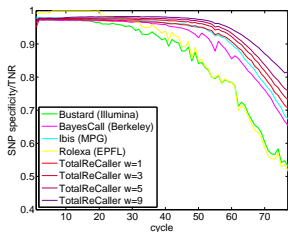


## Phi-X





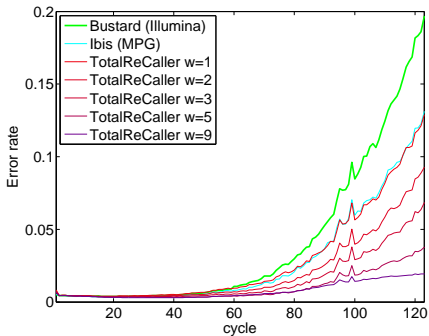
## Phi-X





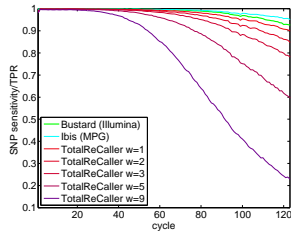
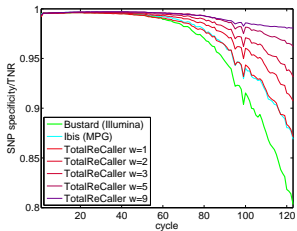


## E.Coli





## E.Coli





## Literatur



M. Burrows and D.J. Wheeler.  
A block-sorting lossless data compression algorithm, 1994.



P. Ferragina and G. Manzini.  
Opportunistic data structures with applications.  
*ANNUAL SYMPOSIUM ON FOUNDATIONS OF COMPUTER SCIENCE*,  
41:390–398, 2000.



M. Kircher, U. Stenzel, and J. Kelso.  
Improved base calling for the Illumina Genome Analyzer using machine learning strategies.  
*Genome Biology*, 10(8):R83, 2009.



W.C. Kao, K. Stevens, and Y.S. Song.  
BayesCall: A model-based base-calling algorithm for high-throughput short-read sequencing.  
*Genome Research*, 19(10):1884, 2009.



Michael L Metzker.  
Emerging technologies in DNA sequencing.  
*Genome research*, 15(12):1767–76, 2005.