

# Nanopore automata

Ian Holmes<sup>1,2,\*</sup>

<sup>1</sup> Lawrence Berkeley National Laboratory, Berkeley, CA, USA

<sup>2</sup> Department of Bioengineering, University of California, Berkeley, CA, USA

## Contents

<b>1</b>	<b>Abstract</b>	<b>2</b>
<b>2</b>	<b>Specification</b>	<b>2</b>
2.1	Parameterization algorithm . . . . .	2
2.2	Reference search algorithm . . . . .	3
2.3	Implementation . . . . .	3
2.4	Evaluation . . . . .	3
<b>3</b>	<b>Methods</b>	<b>3</b>
3.1	Model . . . . .	4
3.2	Baum-Welch algorithm . . . . .	5
3.3	Viterbi algorithm . . . . .	5
<b>4</b>	<b>Results</b>	<b>6</b>
<b>5</b>	<b>Discussion</b>	<b>6</b>
<b>6</b>	<b>Acknowledgments</b>	<b>7</b>
<b>7</b>	<b>Figure Legends</b>	<b>8</b>
<b>8</b>	<b>Appendix</b>	<b>9</b>
8.1	Exponential distribution . . . . .	9
8.2	Gamma distribution . . . . .	9

# 1 Abstract

State machine algorithms for aligning Nanopore reads.

# 2 Specification

Initial goal (Preliminary Results) is simple reusable code for aligning a segmented nanopore read (with segment currents summarized) to a reference sequence.

Longer-term goals (Specific Aims) include

- quasi-hierarchical series of models for processed→raw data (raw, FAST5, FASTQ, FASTA)
- transducer intersection-style models for read-pair alignment
- systematic strategies for approximation/optimization algorithms, climbing the hierarchy (starting with k-mer or FM-index approaches)
- transducer intersection models for aligning reads from different sequencing technologies

## 2.1 Parameterization algorithm

Given the following inputs

- Reference genome (FASTA)
- Segment-called reads (FAST5/HDF5)

Perform the following steps

- Perform Baum-Welch to fit a rich model

Rich model incorporates segment statistics.

## 2.2 Reference search algorithm

Given the following inputs

- Reference genome
- Segment-called reads (FAST5/HDF5)
- Parameterized rich model

Perform the following steps

- Perform Viterbi alignment

## 2.3 Implementation

Libraries etc.

HDF5...

## 2.4 Evaluation

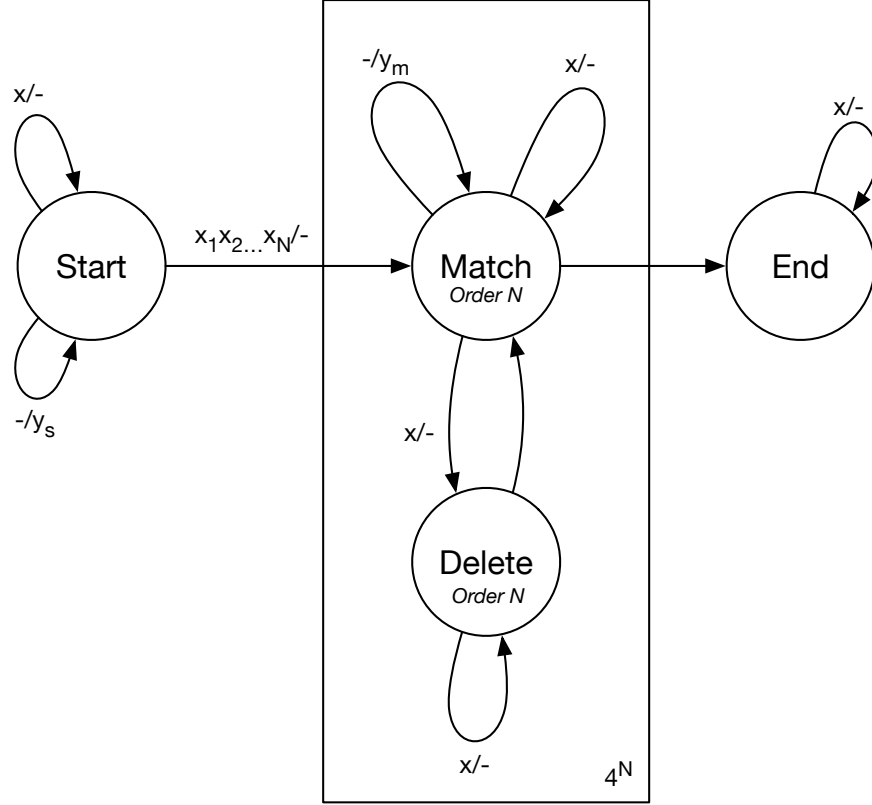
Strategy...

Data sets...

# 3 Methods

Model & inference algorithms.

### 3.1 Model



- Order- $N$  transducer.
- Input: nucleotide
- Output: current levels
- States: **Start**, **Match** $_{x_1 \dots x_N}$ , **Delete** $_{x_1 \dots x_N}$ , **End**
- Transitions
  - **Start**→**Start**: weight  $p^{\text{StartEmit}}$ , emits current  $y_s \sim \text{Normal}(\mu^{\text{Start}}, \tau^{\text{Start}})$
  - **Start**→**Start**: weight 1, absorbs base  $x$
  - **Start**→**Match** $_{x_1 \dots x_N}$ : weight  $1-p^{\text{StartEmit}}$ , absorbs  $N$  bases  $x_1 \dots x_N$

- $\text{Match}_{x_1 \dots x_N} \rightarrow \text{Match}_{x_1 \dots x_N}$ : weight  $p_{x_1 \dots x_N}^{\text{MatchEmit}}$ , emits current  $y_m \sim \text{Normal}(\mu_{x_1 \dots x_N}^{\text{Match}}, \tau_{x_1 \dots x_N}^{\text{Match}})$
- $\text{Match}_{x_1 \dots x_N} \rightarrow \text{Match}_{x_2 \dots x_{N+1}}$ : weight  $(1 - p_{x_1 \dots x_N}^{\text{MatchEmit}})(1 - p^{\text{BeginDelete}})$ , absorbs base  $x_{N+1}$
- $\text{Match}_{x_1 \dots x_N} \rightarrow \text{Delete}_{x_2 \dots x_{N+1}}$ : weight  $(1 - p_{x_1 \dots x_N}^{\text{MatchEmit}})p^{\text{BeginDelete}}$ , absorbs base  $x_{N+1}$
- $\text{Match}_{x_1 \dots x_N} \rightarrow \text{End}$ : weight  $1 - p_{x_1 \dots x_N}^{\text{MatchEmit}}$
- $\text{Delete}_{x_1 \dots x_N} \rightarrow \text{Delete}_{x_2 \dots x_{N+1}}$ : weight  $p^{\text{ExtendDelete}}$ , absorbs base  $x_{N+1}$
- $\text{Delete}_{x_1 \dots x_N} \rightarrow \text{Match}_{x_1 \dots x_N}$ : weight  $1 - p^{\text{ExtendDelete}}$
- $\text{End} \rightarrow \text{End}$ : weight 1, absorbs base  $x$

### 3.2 Baum-Welch algorithm

As usual.

### 3.3 Viterbi algorithm

As usual.

## 4 Results

## 5 Discussion

## **6 Acknowledgments**

## **7 Figure Legends**



## 8 Appendix

### 8.1 Exponential distribution

$$\begin{aligned}
 x &\sim \text{Exponential}(\kappa) \\
 P(x|\kappa) &= \kappa \exp(-\kappa x) \\
 \mathbb{E}[x] &= \kappa^{-1} \\
 \text{Var}[x] &= \kappa^{-2}
 \end{aligned}$$

Rate parameter  $\kappa$ .

### 8.2 Gamma distribution

$$\begin{aligned}
 x &\sim \text{Gamma}(\alpha, \beta) \\
 P(x|\alpha, \beta) &= \frac{x^{\alpha-1} \beta^\alpha \exp(-x\beta)}{\Gamma(\alpha)} \\
 \mathbb{E}[x] &= \alpha/\beta \\
 \text{Var}[x] &= \alpha/\beta^2
 \end{aligned}$$

Shape parameter  $\alpha$ , rate parameter  $\beta$ .  $\Gamma()$  is the gamma function

$$\Gamma(\alpha) = \int_0^\infty z^{\alpha-1} \exp(-z) dz$$

Note  $\Gamma(n) = (n-1)!$  for positive integer  $n$ .

### 8.3 Normal distribution

$$x \sim \text{Normal}(\mu, \tau)$$

Mean  $\mu$ , precision  $\tau$  (precision is reciprocal of variance).

$$P(x|\mu, \tau) = \sqrt{\frac{\tau}{2\pi}} \exp\left(-\frac{\tau}{2}(x - \mu)^2\right)$$