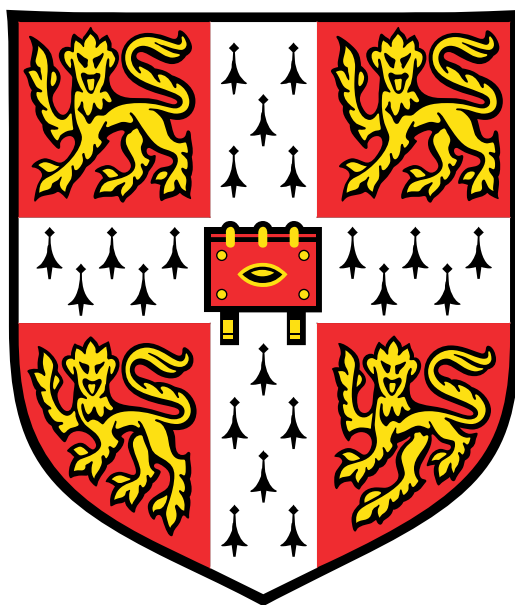


Single molecule mutation detection



Sangjin Lee

Wellcome Trust Sanger Institute
University of Cambridge

This dissertation is submitted for the degree of
Doctor of Philosophy

Downing College

August 2022

I would like to dedicate this thesis to my loving parents ...

Declaration

I hereby declare that except where specific reference is made to the work of others, the contents of this dissertation are original and have not been submitted in whole or in part for consideration for any other degree or qualification in this, or any other university. This dissertation is my own work and contains nothing which is the outcome of work done in collaboration with others, except as specified in the text and Acknowledgements. This dissertation contains fewer than 65,000 words including appendices, bibliography, footnotes, tables and equations and has fewer than 150 figures.

Sangjin Lee
August 2022

Acknowledgements

And I would like to acknowledge ...

Abstract

This is where you write your abstract ...

Table of contents

List of figures	xv
------------------------	-----------

List of tables	xvii
-----------------------	-------------

1 Introduction	1
1.1 Genetic variation and sources of mutations	2
1.2 Germline and somatic mutation	2
1.2.1 Germline mutations	2
1.2.2 Somatic mutations	2
1.2.3 Mutational signatures and mutational processes	2
1.3 Reference genomes	2
1.3.1 De novo assembly	2
1.3.2 Short-read sequencing	2
1.3.3 Long-read sequencing	2
1.3.4 Haplotype tagging	2
1.3.5 High-throughput chromatin conformation capture sequencing	2
1.4 Resequencing	2
1.4.1 Germline mutation detection	2
1.4.2 Somatic mutation detection	2
1.4.3 Somatic mutation detection in cancer	2
1.4.4 Somatic mutation detection in normal tissues	2
1.5 Darwin Tree of Life project	2
1.6 Sperm sequencing for meiotic recombination product investigation	2
1.7 Bloom syndrome patient sample sequencing for mitotic gene conversion detection	2
1.8 Overview and objectives	2

2	Single molecule somatic single-base substitution detection	3
2.1	Introduction	4
2.2	Materials & Methods	4
2.2.1	Samples	4
2.2.2	Mutational signatures	4
2.2.3	Pacific Biosciences Circular Consensus Sequencing	4
2.2.4	Single molecule somatic mutation calling	4
2.2.5	Methods for single molecule somatic mutation calling	4
2.2.6	Mutation calling	4
2.2.7	Hard filters	4
2.2.8	Haplotype Phasing	4
2.3	Benchmarks	4
2.3.1	Sensitivity and Specificity, F1-statistics	4
2.3.2	Receiver-operating characteristics	4
2.4	Results	4
2.4.1	4
2.4.2	4
2.5	Discussion	4
2.5.1	Liquid Biopsy	4
2.5.2	False positive substitutions	4
2.5.3	DeepConsensus	4
2.5.4	Environmental mutagenesis	4
2.5.5	Single molecule structural variation detection	4
2.5.6	Single molecule somatic mutation detection in non-human species	4
3	Germline and somatic mutational processes across the Tree of Life	5
3.1	Introduction	5
3.2	Materials & Methods	5
3.3	Results	5
3.3.1	Germline and somatic mutational processes	5
3.4	Discussion	5
3.4.1	Samples without somatic mutations	5
3.4.2	Somatic theory of aging	5
3.4.3	Life cycle of Insects	5
3.4.4	Environmental mutagenesis	5

4	Meiotic recombination	7
4.1	Introduction	7
4.1.1	Meiotic recombination	7
4.1.2	Haplotype Map	7
4.1.3	Methods to study meiotic recombinant products	7
4.2	Material & Methods	7
4.3	Results	7
4.4	Discussion	7
5	Mitotic Loss of Heterozygosity as an Oncogenic Mechanism	9
5.1	Introduction	9
5.2	Materials & Methods	9
5.3	Discussion	9
Appendix A	How to install L^AT_EX	11
Appendix B	Installing the CUED class file	15

List of figures

List of tables

Chapter 1

Introduction

1.1 Genetic variation and sources of mutations

1.2 Germline and somatic mutation

1.2.1 Germline mutations

1.2.2 Somatic mutations

1.2.3 Mutational signatures and mutational processes

1.3 Reference genomes

1.3.1 De novo assembly

1.3.2 Short-read sequencing

1.3.3 Long-read sequencing

1.3.4 Haplotype tagging

1.3.5 High-throughput chromatin conformation capture sequencing

1.4 Resequencing

1.4.1 Germline mutation detection

1.4.2 Somatic mutation detection

1.4.3 Somatic mutation detection in cancer

1.4.4 Somatic mutation detection in normal tissues

Single-cell expansion and sequencing

Chapter 2

Single molecule somatic single-base substitution detection

2.1 Introduction

2.2 Materials & Methods

2.2.1 Samples

2.2.2 Mutational signatures

2.2.3 Pacific Biosciences Circular Consensus Sequencing

2.2.4 Single molecule somatic mutation calling

2.2.5 Methods for single molecule somatic mutation calling

2.2.6 Mutation calling

2.2.7 Hard filters

2.2.8 Haplotype Phasing

2.3 Benchmarks

2.3.1 Sensitivity and Specificity, F1-statistics

2.3.2 Receiver-operating characteristics

2.4 Results

2.4.1

2.4.2

Chapter 3

Germline and somatic mutational processes across the Tree of Life

3.1 Introduction

3.2 Materials & Methods

3.3 Results

3.3.1 Germline and somatic mutational processes

3.4 Discussion

3.4.1 Samples without somatic mutations

3.4.2 Somatic theory of aging

3.4.3 Life cycle of Insects

3.4.4 Environmental mutagenesis

Chapter 4

Meiotic recombination

4.1 Introduction

4.1.1 Meiotic recombination

4.1.2 Haplotype Map

4.1.3 Methods to study meiotic recombinant products

Trio-sequencing

4.2 Material & Methods

4.3 Results

4.4 Discussion

Chapter 5

Mitotic Loss of Heterozygosity as an Oncogenic Mechanism

5.1 Introduction

5.2 Materials & Methods

5.3 Discussion

Appendix A

How to install L^AT_EX

Windows OS

TeXLive package - full version

1. Download the TeXLive ISO (2.2GB) from
<https://www.tug.org/texlive/>
2. Download WinCDEmu (if you don't have a virtual drive) from
<http://wincdemu.sysprogs.org/download/>
3. To install Windows CD Emulator follow the instructions at
<http://wincdemu.sysprogs.org/tutorials/install/>
4. Right click the iso and mount it using the WinCDEmu as shown in
<http://wincdemu.sysprogs.org/tutorials/mount/>
5. Open your virtual drive and run setup.pl

or

Basic MikTeX - T_EX distribution

1. Download Basic-MiK_TE_X(32bit or 64bit) from
<http://miktex.org/download>
2. Run the installer
3. To add a new package go to Start » All Programs » MikTeX » Maintenance (Admin)
and choose Package Manager

4. Select or search for packages to install

TexStudio - T_EX editor

1. Download TexStudio from
<http://texstudio.sourceforge.net/#downloads>
2. Run the installer

Mac OS X

MacTeX - T_EX distribution

1. Download the file from
<https://www.tug.org/mactex/>
2. Extract and double click to run the installer. It does the entire configuration, sit back and relax.

TexStudio - T_EX editor

1. Download TexStudio from
<http://texstudio.sourceforge.net/#downloads>
2. Extract and Start

Unix/Linux

TeXLive - T_EX distribution

Getting the distribution:

1. TeXLive can be downloaded from
<http://www.tug.org/texlive/acquire-netinstall.html>.
2. TeXLive is provided by most operating system you can use (rpm,apt-get or yum) to get TeXLive distributions

Installation

1. Mount the ISO file in the mnt directory

```
mount -t iso9660 -o ro,loop,noauto /your/texlive####.iso /mnt
```

2. Install wget on your OS (use rpm, apt-get or yum install)
3. Run the installer script install-tl.

```
cd /your/download/directory
./install-tl
```

4. Enter command 'i' for installation
5. Post-Installation configuration:
<http://www.tug.org/texlive/doc/texlive-en/texlive-en.html#x1-320003.4.1>
6. Set the path for the directory of TexLive binaries in your .bashrc file

For 32bit OS

For Bourne-compatible shells such as bash, and using Intel x86 GNU/Linux and a default directory setup as an example, the file to edit might be

```
edit ~/.bashrc file and add following lines
PATH=/usr/local/texlive/2011/bin/i386-linux:$PATH;
export PATH
MANPATH=/usr/local/texlive/2011/texmf/doc/man:$MANPATH;
export MANPATH
INFOPATH=/usr/local/texlive/2011/texmf/doc/info:$INFOPATH;
export INFOPATH
```

For 64bit OS

```
edit ~/.bashrc file and add following lines
PATH=/usr/local/texlive/2011/bin/x86_64-linux:$PATH;
export PATH
MANPATH=/usr/local/texlive/2011/texmf/doc/man:$MANPATH;
export MANPATH
```

```
INFOPATH=/usr/local/texlive/2011/texmf/doc/info:$INFOPATH;  
export INFOPATH
```

Fedora/RedHat/CentOS:

```
sudo yum install texlive  
sudo yum install psutils
```

SUSE:

```
sudo zypper install texlive
```

Debian/Ubuntu:

```
sudo apt-get install texlive texlive-latex-extra  
sudo apt-get install psutils
```


Appendix B

Installing the CUED class file

\LaTeX .cls files can be accessed system-wide when they are placed in the $\langle\text{texmf}\rangle/\text{tex}/\text{latex}$ directory, where $\langle\text{texmf}\rangle$ is the root directory of the user's \TeX installation. On systems that have a local texmf tree ($\langle\text{texmflocal}\rangle$), which may be named “ texmf-local ” or “ localtexmf ”, it may be advisable to install packages in $\langle\text{texmflocal}\rangle$, rather than $\langle\text{texmf}\rangle$ as the contents of the former, unlike that of the latter, are preserved after the \LaTeX system is reinstalled and/or upgraded.

It is recommended that the user create a subdirectory $\langle\text{texmf}\rangle/\text{tex}/\text{latex}/\text{CUED}$ for all CUED related \LaTeX class and package files. On some \LaTeX systems, the directory look-up tables will need to be refreshed after making additions or deletions to the system files. For \TeX Live systems this is accomplished via executing “ texhash ” as root. MikTeX users can run “ initexmf -u ” to accomplish the same thing.

Users not willing or able to install the files system-wide can install them in their personal directories, but will then have to provide the path (full or relative) in addition to the filename when referring to them in \LaTeX .

