# Transcriptome and isoform

# reconstruction with long reads

**Adam Ameur, PhD**
**NGI, Uppsala Genome Center**
**adam.ameur@igp.uu.se**

RNA-seq workshop 2014

# National Genomics Infrastructure



NGI Stockholm



NGI Uppsala

**NGI staff**: 60 -70 FTE, including head of facility, lab research engineers, bioinformaticians, IT-experts, project coordinators.

**UPPMAX/UPPNEX**: Uppsala multidisciplinary center for advanced computational science, UPPNEX: UPPmax NEXt generation sequencing Cluster & Storage.

Enabler for Life Sciences

Karolinska Institutet

KTH VETENSKAP OCH KONST

Stockholms universitet

UPPSALA UNIVERSITET

SciLifeLab

# DNA sequencing at all scales



One of the most well-equipped NGS sites in Europe!

| | |
|---|---|
| 10 | **Illumina HiSeq Xten** |
| 17 | **Illumina HiSeq 2000/2500** |
| 3 | **Illumina MiSeq** |
| 1 | **Illumina NextSeq** |
| 2 | **Life Technologies Ion Torrent** |
| 6 | **Life Technologies Ion Proton** |
| 2 | **Pacific Biosciences RSII** |
| 2 | **Sanger ABI3730** |
| 1 | **Argus Whole Genome Map. Syst.** |
| 1 | **Oxford Nanopore MinIon** |

Enabler for Life Sciences

SciLifeLab

# RNA-sequencing

# with short reads

Karolinska Institutet · KTH · Stockholms universitet · Uppsala Universitet

SciLifeLab

# **RNA-seq:** standard procedure



fragmen-
tation

RT

mRNA

RT

fragmen-
tation

sequence library

short sequence reads
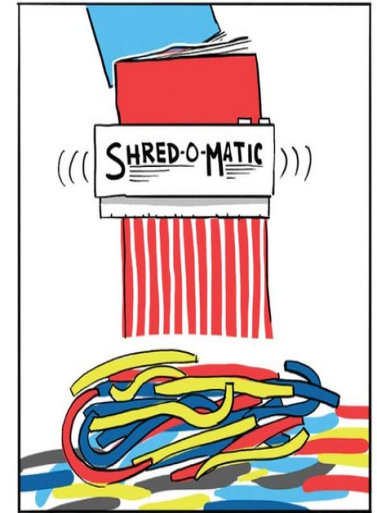
# **RNA-seq:** the main question

What to do with this?

# RNA-seq: analysis

# Complicating factor: alternative splicing

# RNA-seq: problem with short reads

# RNA-sequencing
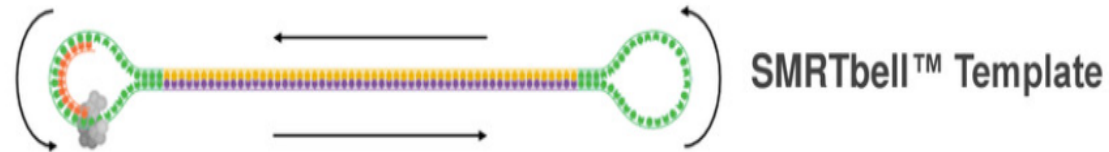
# with very long reads!!!

# Pacific Biosciences RS II

- Pacific Biosciences
    - Single molecule sequencing
    - Very long read lengths (up to 30 kb)
    - Rapid sequencing
    - Can detect base modifications (e.g. methylation)
    - Relatively low throughput

Pacific Biosciences RSII

# PacBio – Sequencing Template



SMRTbell™ Template

## Polymerase Read

### Definition:
- Sequence of nucleotides incorporated by polymerase while reading a template
- Includes adapters
- Often called "read"
- Includes adapters
- 1 molecule, 1 pol. read

### Uses:
- QC of instrument run
- Benchmarking

## Subread

### Definition:
- Single pass of template
- Adapters removed
- 1 molecule, >=1 subread

### Unique data:
- Kinetic measurements
- Rich QVs

### Uses:
- Applications

## Read (of Insert)

### Definition:
- Represents highest-quality single-sequence for an insert, regardless of number of passes
- Generalizes CCS for <2 passes & RQ <0.9
- 1 or more passes
- 1 molecule, 1 read

### Uses:
- Library QC
- Applications

# PacBio output

- ## PacBio throughput
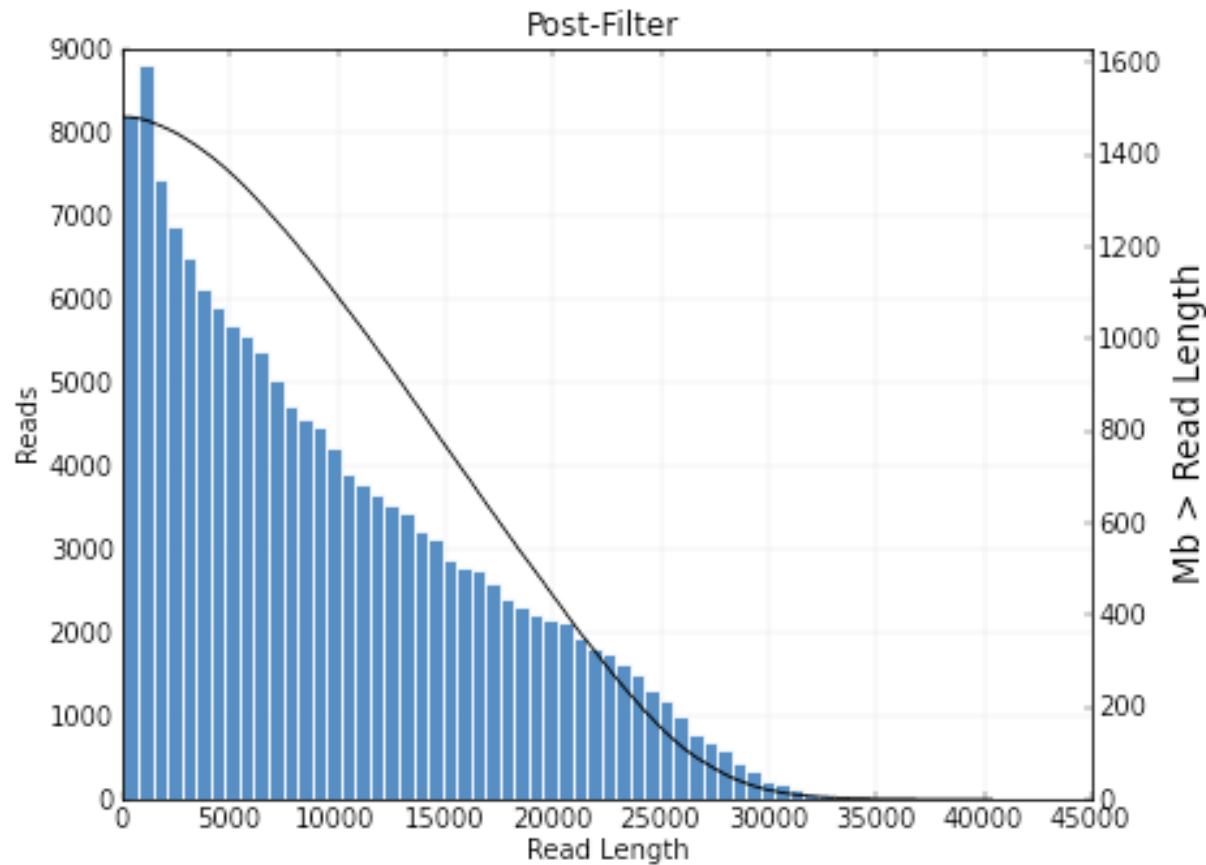  ~ 500Mb-1Gb/SMRT cell

  ~1 bacterial genome

  ~1 bacterial transcriptome

  1 human genome = 150 SMRT cells

- ## PacBio read lengths: 500bp-30kb

# PacBio – Current read lengths

- >10kb average read lengths! (run from April 2014)

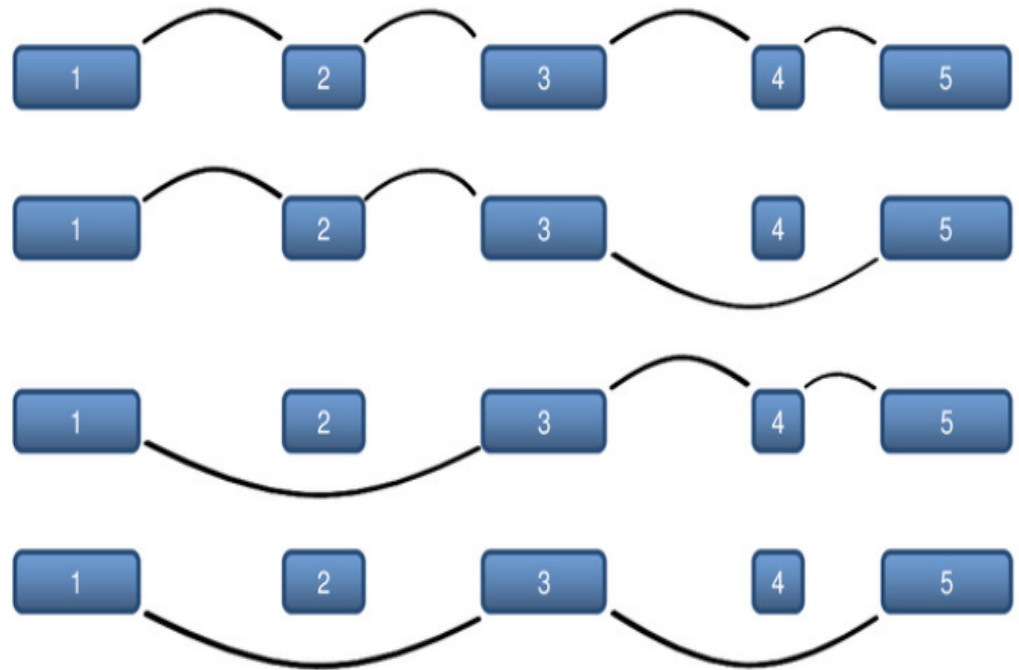# Iso-Seq: Full length RNA-seq on PacBio!

- ## Single molecule sequencing
  - One read – one transcript

- ## Transcript in full length
  - No assembly required

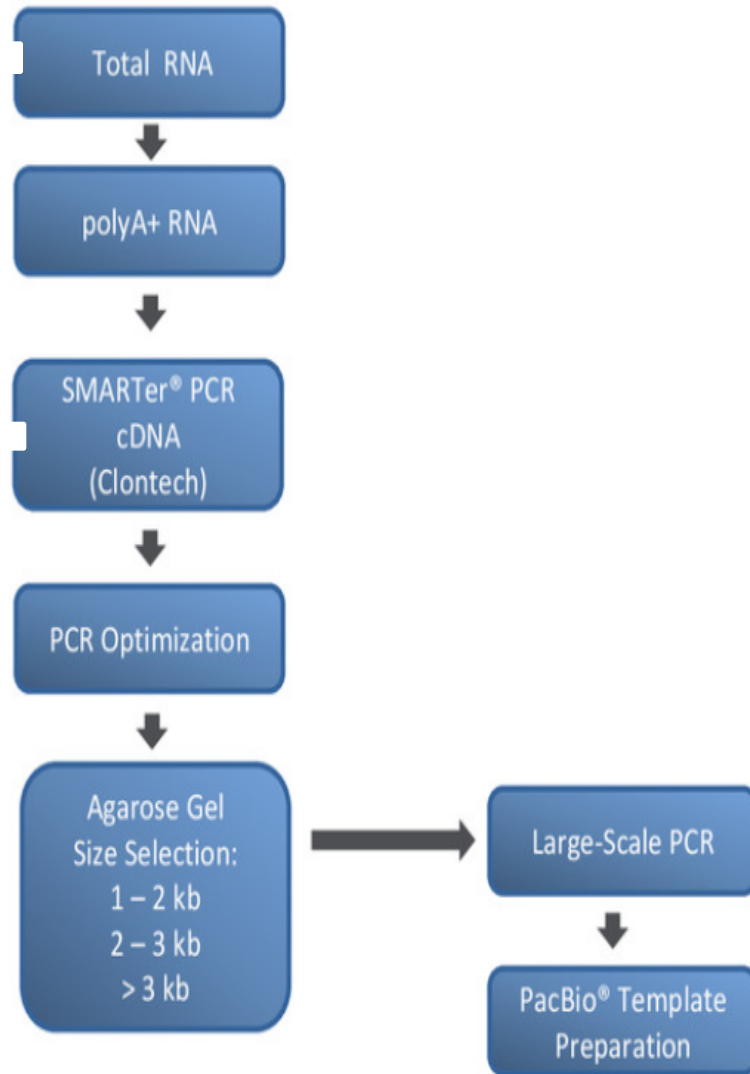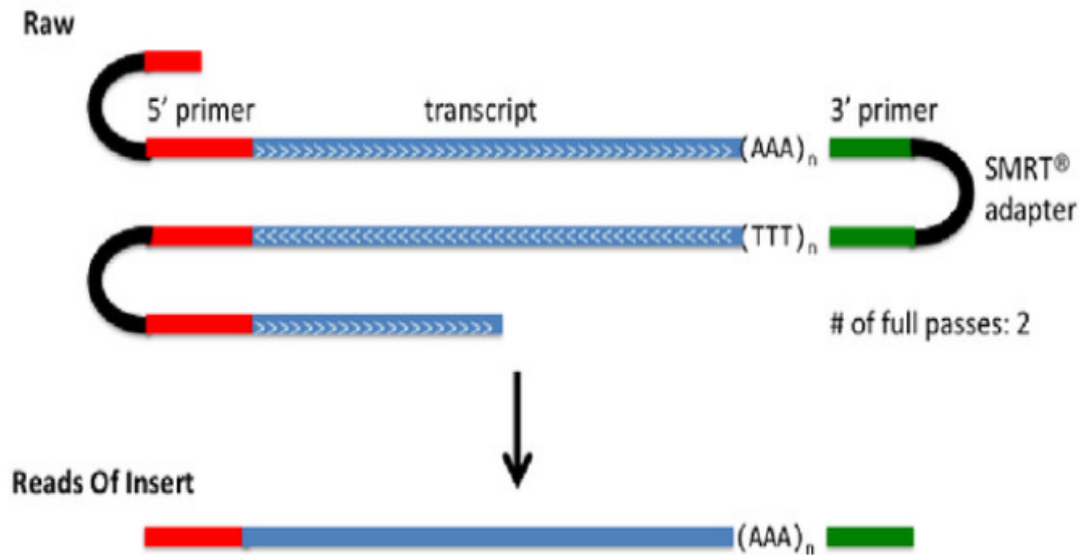- ## No systematic bias
  - CG-rich, AT-rich, tandem repeats

# PacBio Iso-Seq - library preparation
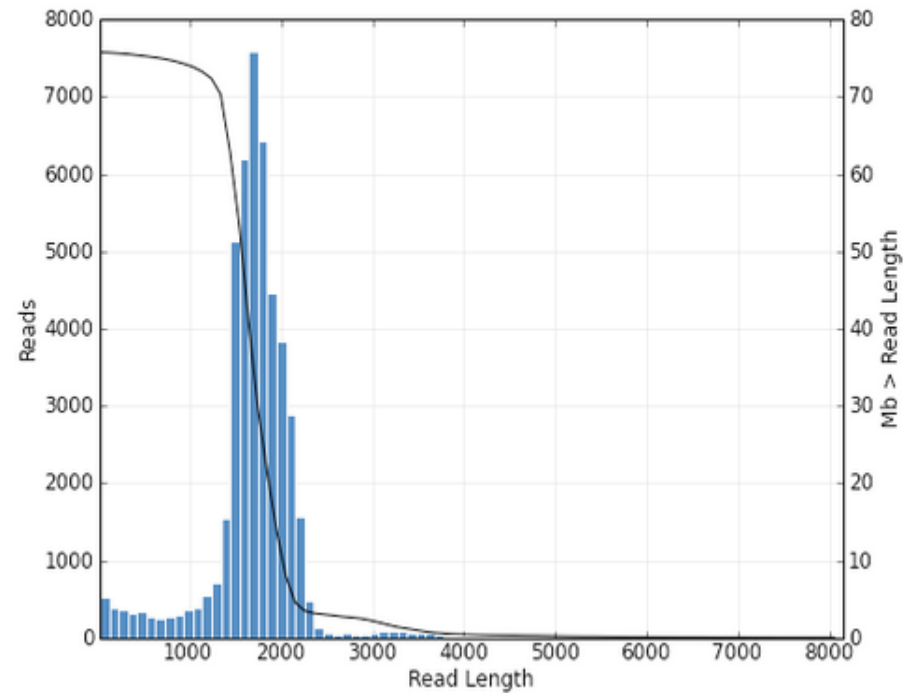
# PacBio Iso-Seq – reads of insert



Full-Length = 5' primer seen, polyA tail seen, 3' primer seen

- Identify and remove primers and polyA/T tail
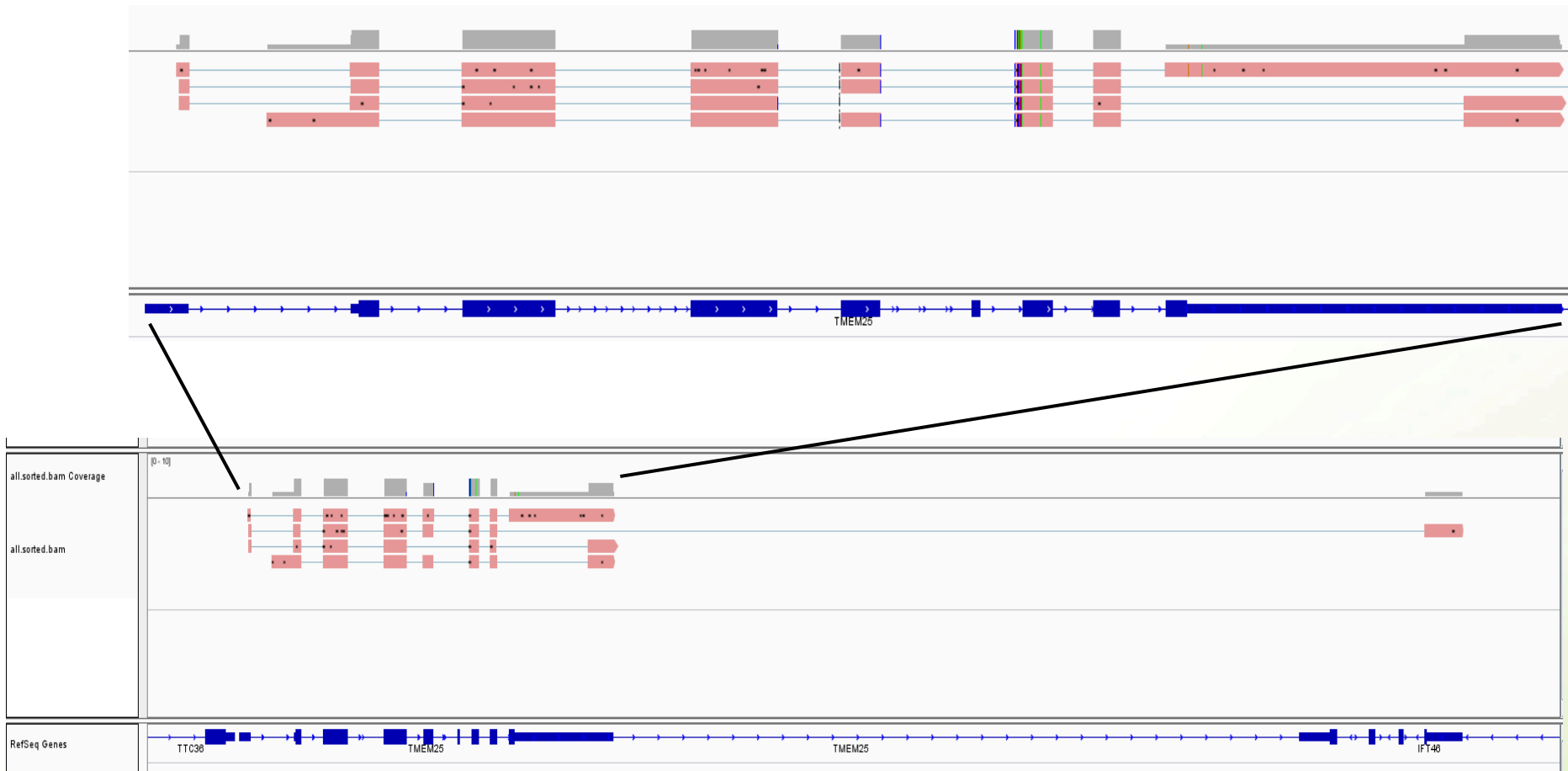- Identify read strandedness

# PacBio Iso-Seq: ROI of 1kb lib (2 cells)

| | |
|---|---|
| Read Bases of Insert | 78,108,189 |
| Mean Read Length of Insert | 1,687 |
| Mean Number of Passes | **8.0** |
| Number of full-length non-chimeric reads | **35,467** |
| Average full-length non-chimeric read length | 1,679 |



Read Length Of Insert

# PacBio Iso-Seq: examples



GeneCards Summary for TMEM25 Gene:
**TMEM25** is a protein-coding gene. Diseases associated with TMEM25 include breast cancer.

# PacBio Iso-Seq: examples (from PacBio)

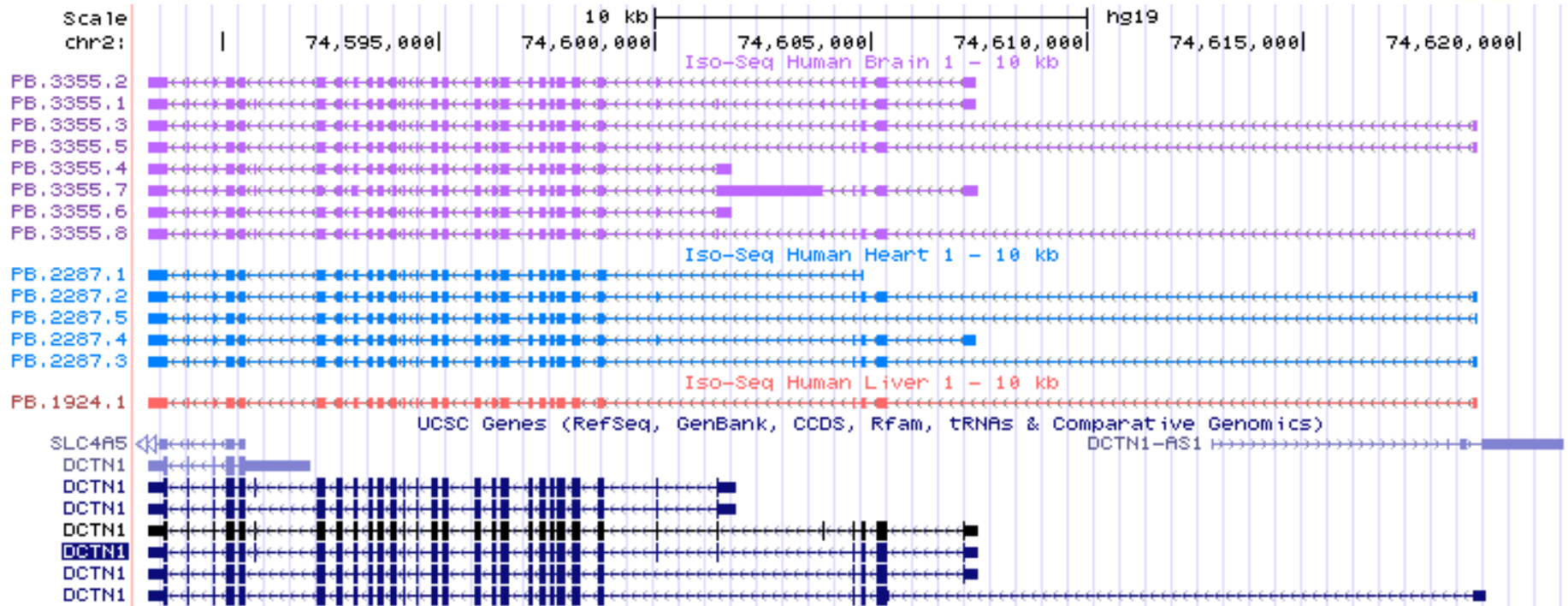| Tissue | Size Selection | FL Reads | Average FL Readlength | Number of Unique FL Transcripts | Number of Gene Loci | Max Transcript Length |
|---|---|---|---|---|---|---|
| Brain | 1 - 2 kb | 159792 | 1785 | 10289 | 6356 | 8823 |
| | 2 - 3 kb | 165942 | 2794 | | | |
| | 3 - 6 kb | 118568 | 4104 | | | |
| | 5 - 10 kb | 59607 | 6490 | | | |
| Heart | 1 - 2 kb | 134462 | 1629 | 6896 | 4352 | 8528 |
| | 2 - 3 kb | 89472 | 2910 | | | |
| | 3 - 6 kb | 126927 | 4027 | | | |
| | 5 - 10 kb | 43486 | 6323 | | | |
| Liver | 1 - 2 kb | 197772 | 1725 | 6124 | 3497 | 4754 |
| | 2 - 3 kb | 157531 | 2605 | | | |
| | 3 - 6 kb | 130438 | 3876 | | | |

*http://blog.pacificbiosciences.com/2014/10/data-release-whole-human-transcriptome.html*
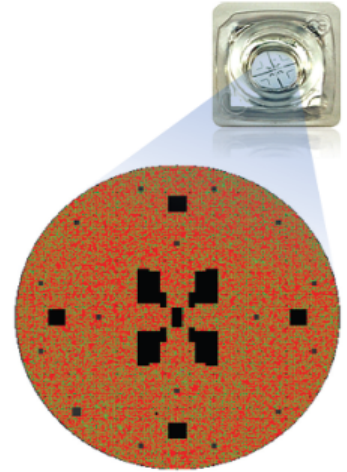
# PacBio Iso-Seq: examples (from PacBio)

| Tissue | Size Selection | FL Reads | Average FL Readlength | Number of Unique FL Transcripts | Number of Gene Loci | Max Transcript Length |
|--------|----------------|----------|-----------------------|--------------------------------|---------------------|-----------------------|
| Brain  | 1 - 2 kb       | 159792   | 1785                  | 10289                          | 6356                | 8823                  |
|        | 2 - 3 kb       | 165942   | 2794                  |                                |                     |                       |
|        | 3 - 6 kb       | 118568   | 4104                  |                                |                     |                       |
|        | 5 - 10 kb      | 59607    | 6490                  |                                |                     |                       |
| Heart  | 1 - 2 kb       | 134462   | 1629                  | 6896                           | 4352                | 8528                  |
|        | 2 - 3 kb       | 89472    | 2910                  |                                |                     |                       |
|        | 3 - 6 kb       | 126927   | 4027                  |                                |                     |                       |
|        | 5 - 10 kb      | 43486    | 6323                  |                                |                     |                       |
| Liver  | 1 - 2 kb       | 197772   | 1725                  | 6124                           | 3497                | 4754                  |
|        | 2 - 3 kb       | 157531   | 2605                  |                                |                     |                       |
|        | 3 - 6 kb       | 130438   | 3876                  |                                |                     |                       |

*http://blog.pacificbiosciences.com/2014/10/data-release-whole-human-transcriptome.html*

# PacBio Iso-Seq: examples (from PacBio)



*http://blog.pacificbiosciences.com/2014/10/data-release-whole-human-transcriptome.html*

# PacBio Iso-Seq: experimental design

So how many SMRT cells do I need?

Approximate scope guidance:
- **1 SMRT Cell**: targeted, gene-specific isoform characterization
- **1-8 SMRT Cells**: get a high-level overview of the transcriptome and isoforms of abundant transcripts
- **8-50 SMRT Cells**: get a detailed look at most transcripts and their isoforms
- **>50 SMRT Cells**: get a very thorough look at transcriptome with rare transcripts and rare isoforms or intermediates

! Depends strongly on transcriptome complexity of the organism being studied !

# Clinical project: Chronic Myeloid Leukemia

- BCR-ABL1 fusion protein – a CML drug target



The BCR-ABL1 fusion protein can acquire resistance mutations following drug treatment

www.cambridgemedicine.org/article/doi/10.7244/cmj-1355057881

# BCR-ABL1 workflow – PacBio Sequencing

BMC Cancer

**RESEARCH ARTICLE**                                                        **Open Access**

## Clonal distribution of *BCR-ABL1* mutations and splice isoforms by single-molecule long-read RNA sequencing

Lucia Cavelier[1*†], Adam Ameur[1†], Susana Häggqvist[1], Ida Höijer[1], Nicola Cahill[1], Ulla Olsson-Strömberg[2] and Monica Hermanson[1]

**Abstract**

**Background:** The evolution of mutations in the *BCR-ABL1* fusion gene transcript renders CML patients resistant to tyrosine kinase inhibitor (TKI) based therapy. Thus screening for *BCR-ABL1* mutations is recommended particularly in patients experiencing poor response to treatment. Herein we describe a novel approach for the detection and surveillance of *BCR-ABL1* mutations in CML patients.

**Methods:** To detect mutations in the *BCR-ABL1* transcript we developed an assay based on the Pacific Biosciences (PacBio) sequencing technology, which allows for single-molecule long-read sequencing of *BCR-ABL1* fusion transcript molecules. Samples from six patients with poor response to therapy were analyzed both at diagnosis and follow-up. cDNA was generated from total RNA and a 1,6 kb fragment encompassing the *BCR-ABL1* transcript was amplified using long range PCR. To estimate the sensitivity of the assay, a serial dilution experiment was performed.

**Results:** Over 10,000 full-length *BCR-ABL1* sequences were obtained for all samples studied. Through the serial dilution analysis, mutations in CML patient samples could be detected down to a level of at least 1%. Notably, the assay was determined to be sufficiently sensitive even in patients harboring a low abundance of *BCR-ABL1* levels. The PacBio sequencing successfully identified all mutations seen by standard methods. Importantly, we identified several mutations that escaped detection by the clinical routine analysis. Resistance mutations were found in all but one of the patients. Due to the long reads afforded by PacBio sequencing, compound mutations present in the same molecule were readily distinguished from independent alterations arising in different molecules. Moreover, several transcript isoforms of the *BCR-ABL1* transcript were identified in two of the CML patients. Finally, our assay allowed for a quick turn around time allowing samples to be reported upon within 2 days.

**Conclusions:** In summary the PacBio sequencing assay can be applied to detect *BCR-ABL1* resistance mutations in both diagnostic and follow-up CML patient samples using a simple protocol applicable to routine diagnosis. The method besides its sensitivity, gives a complete view of the clonal distribution of mutations, which is of importance when making therapy decisions.

# BCR-ABL1 mutations at diagnosis

PacBio sequencing generates ~10 000X coverage!



Sample from time of diagnosis:

# BCR-ABL1 mutations in follow-up sample



Sample 6 months later

Mutations acquired in fusion transcript.
Might require treatment with alternative drug.
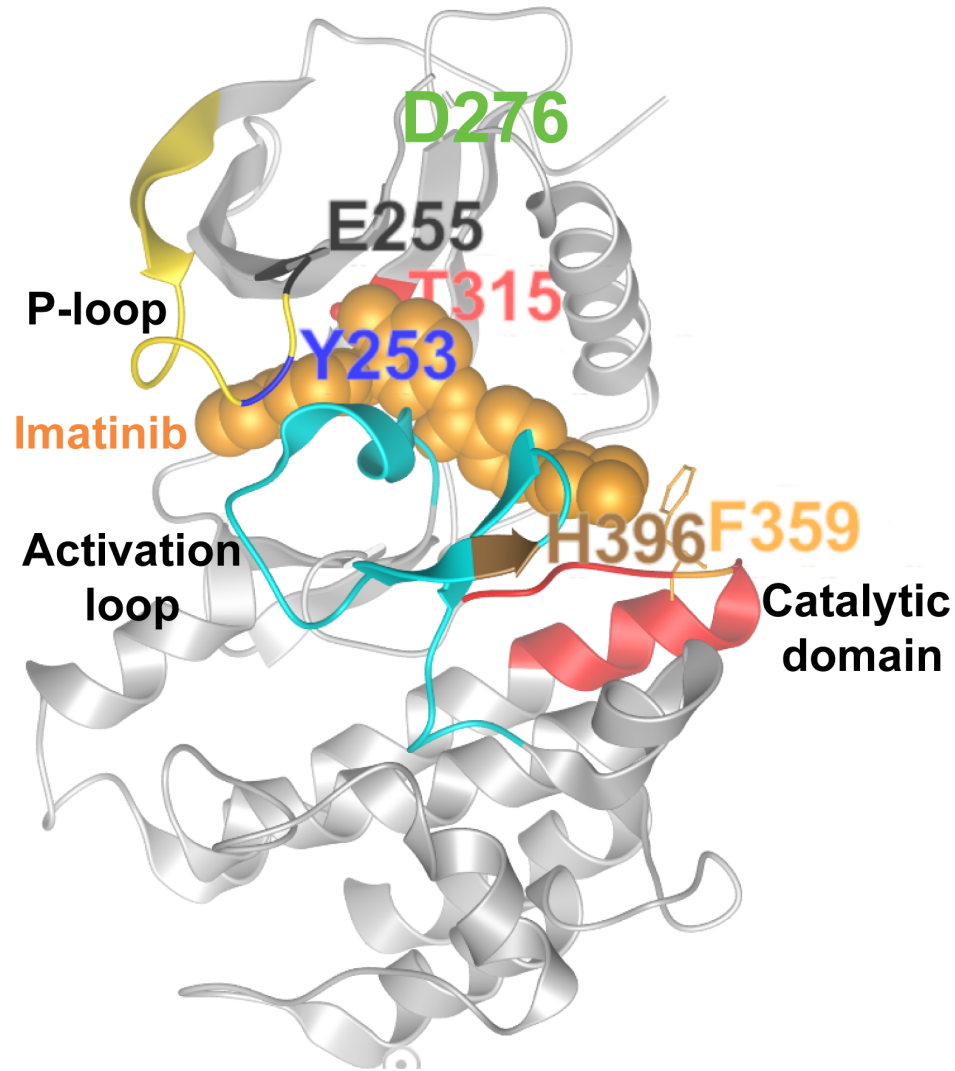
# BCR-ABL1 dilution series results

- Mutations down to 1% detected!

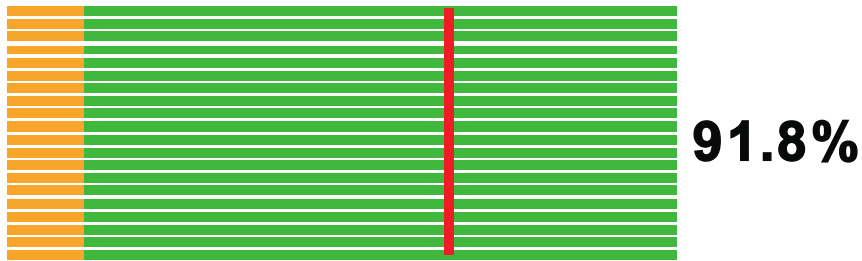# Summary of mutations in 5 CML patients

# Mutations mapped to protein structure
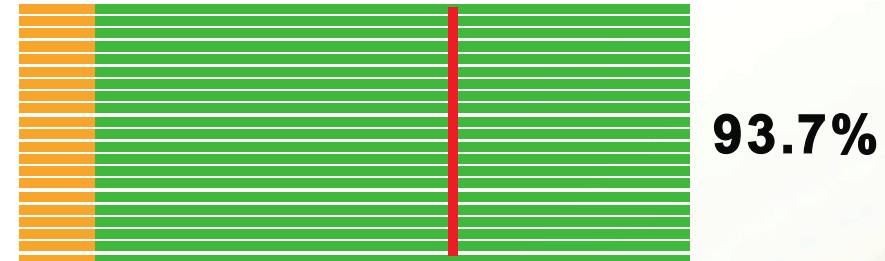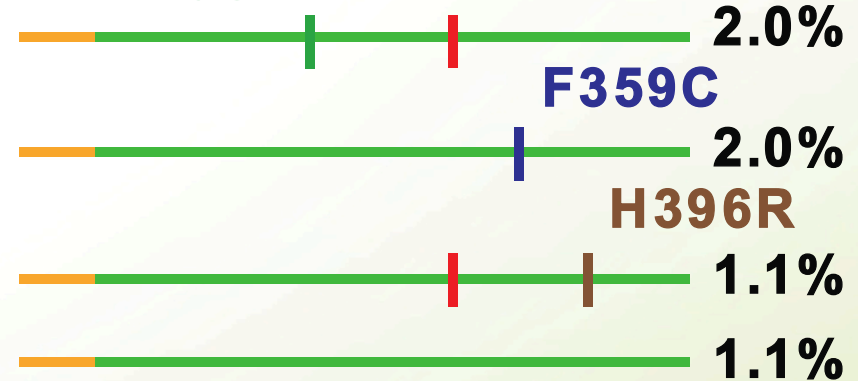
# BCR-ABL1 - Compound mutations

# Future bioinformatics challenge

- How to find mutations within isoforms???

# Conclusions and next steps

- Sensitive method for *BCR-ABL1* analysis!
  - Also for compound mutations and isoforms

- Method now used in clinical routine!
  - Patient samples coming to the clinic over a few months
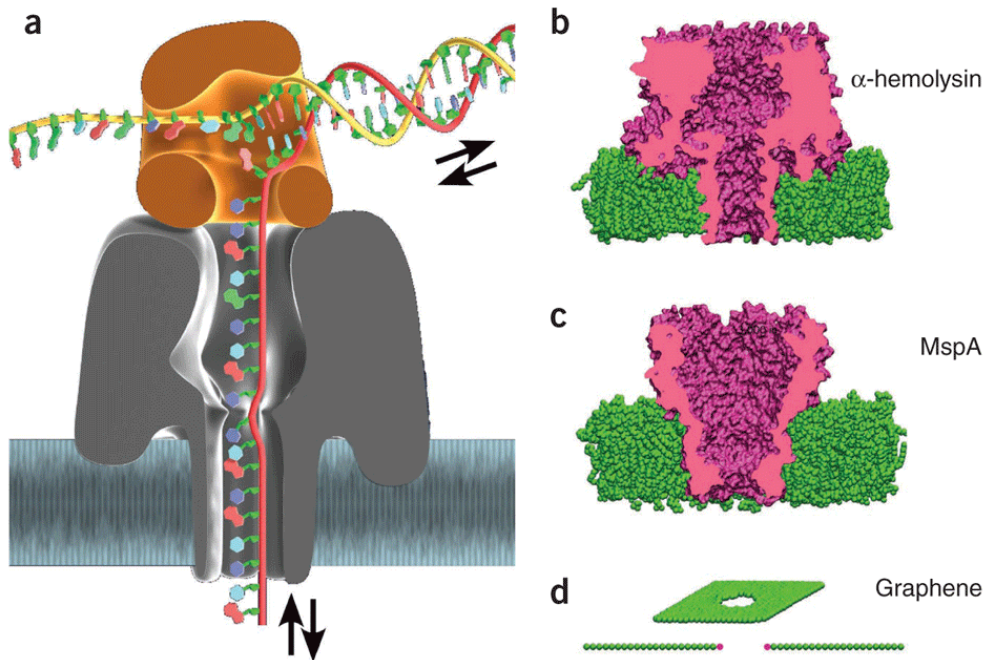  - Response time limit: 2 weeks

**Sequel** - New PacBio instrument with higher throughput!



**7x more data per SMRT cell!**

# News and future directions (2)

Nanopore technology - for direct RNA sequencing?



**PromethION**

Enables detection of modified RNA bases??