

Zhiying Zhao^{1*}, Yu-Chih Tsai², Alicia Clum¹, Katherine Munson¹, Chris Daum¹, Stephen W. Turner², Jonas Korlach², Len A. Pennacchio¹, Feng Chen¹
1. Department of Energy, Joint Genome Institute, 2800 Mitchell Drive, Walnut Creek, CA, 94598. 2. Pacific Biosciences, 1380 Willow Rd., Menlo Park, CA, 94025

INTRODUCTION

The assembly and analysis of microbial species on earth remains a largely unexplored area of life. This is partially due to their inability to be cultured but also based on the large historic cost of drafting and finishing individual microbial species genomes.

The single-molecule real-time (SMRT™) sequencing platform developed by Pacific Biosciences (PacBio) offers several benefits including Single Molecule real-time analysis, longer read length at fast speed, low sequencing redundancy and bias. Thus, it was used at JGI as a quick-turnaround and cost-effective solution for finishing microbial genomes.

Construction of PacBio library by traditional protocol still requires micrograms of genomic DNA. In many cases, getting high quantity of genomic DNA remains as a major challenge. Recently, PacBio developed a more efficient library construction method using terminal deoxynucleotidyl transferase (TdT), which makes it possible to obtain sufficient sequencing data for assembly from significantly smaller amount of genomic DNA. We have tested and validated this newly developed method. Preliminary analysis results suggested that this technology can be used for microbial genome assembly with PacBio only data.

LIBRARY CONSTRUCTION

Ten bacterial samples (various GC% and genome size) are selected for validation. The library creation process begins with fragmenting genomic DNA (100-200ng) to 10kb using Covaris G-tube, followed by damage repair, quick exonuclease treatment and PolyA tailing. Ampure SPRI beads are used throughout library preparation process to select and purify sample DNA. The total preparation time is shorter than standard SMRTbell library construction processes. The library could be constructed within 4 hours.

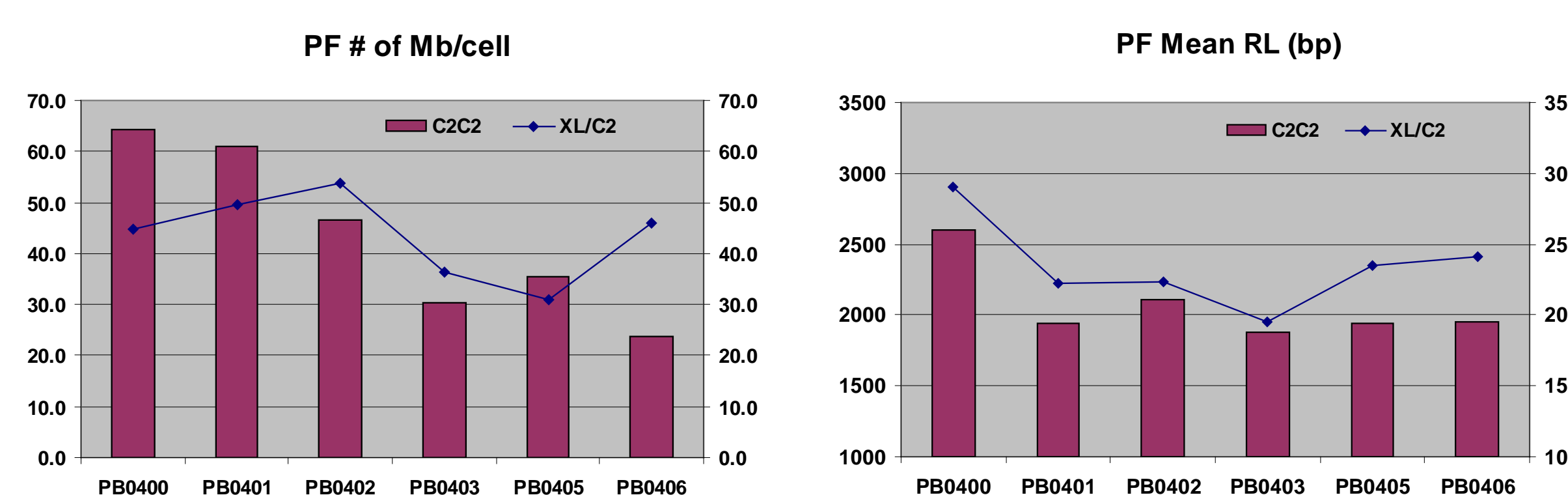
Sample Info			Input	Library Info	
Sample Name	GC%	Genome Size (MB)	ng	ng	Yield %
Tolomonas sp. BRL6-1	47	4.1	100	81	81%
Gillisia sp. JM1	34	5.4	100	72	72%
Teredinibacter sp. strain 991H.S.0a.06	50	7.2	200	160	80%
Geopsychrobacter electrodiphilus DSM 16401	53	5.0	200	136	68%
Hippea medeae KM1	43	4.7	181	123	68%
Desulfospira joergensenii DSM 10085	50	6.3	70	45	65%
Streptomyces sp. WmmB714	72	6.6	200	152	76%
Nocardia sp. BMG111209	69	9.1	200	127	63%
Nocardia sp. BMG51109	68	8.8	116	94	81%
Meiothermus ruber DSM 1279	63	3.1	100	73	73%

SEQUENCING RUN AND RESULTS

Sequencing run was done with Magbead loading, stage start, and 120min movies on either V2 or V3 chips, targeting 100x coverage per genome.

Sample Name	GC%	Genome Size (MB)	Sequencing Chemistry (# of Cells)		Sequencing Results			
			C2/C2	XL/C2	PF # of Reads/cel	PF Mean RL (bp)	PF # of Mb/cell	PF Mean RQ
Tolomonas sp. BRL6-1	47	4.1	12	4	10225	2950	29.9	82.5%
Gillisia sp. JM1	34	5.4	12	4	11714	2848	32.9	81.2%
Teredinibacter sp. strain 991H.S.0a.06	50	7.2	6	6	20000	2755	54.4	83.1%
Geopsychrobacter electrodiphilus DSM 16401	53	5.0	4	5	26325	2101	54.7	83.0%
Hippea medeae KM1	43	4.7	6	2	22597	2156	48.4	79.8%
Desulfospira joergensenii DSM 10085	50	6.3	6	5	17268	1911	33.0	80.1%
Nocardia sp. BMG111209	69	9.1	8	8	15729	2146	33.3	82.3%
Nocardia sp. BMG51109	68	8.8	4	8	16762	2263	38.6	82.3%
Meiothermus ruber DSM 1279	63	3.1	6	0	29513	1899	55.7	83.3%

Differences between sequencing chemistries: XL/C2 tends to give longer read length. There is no clear trend for per cell output in terms number of bases.

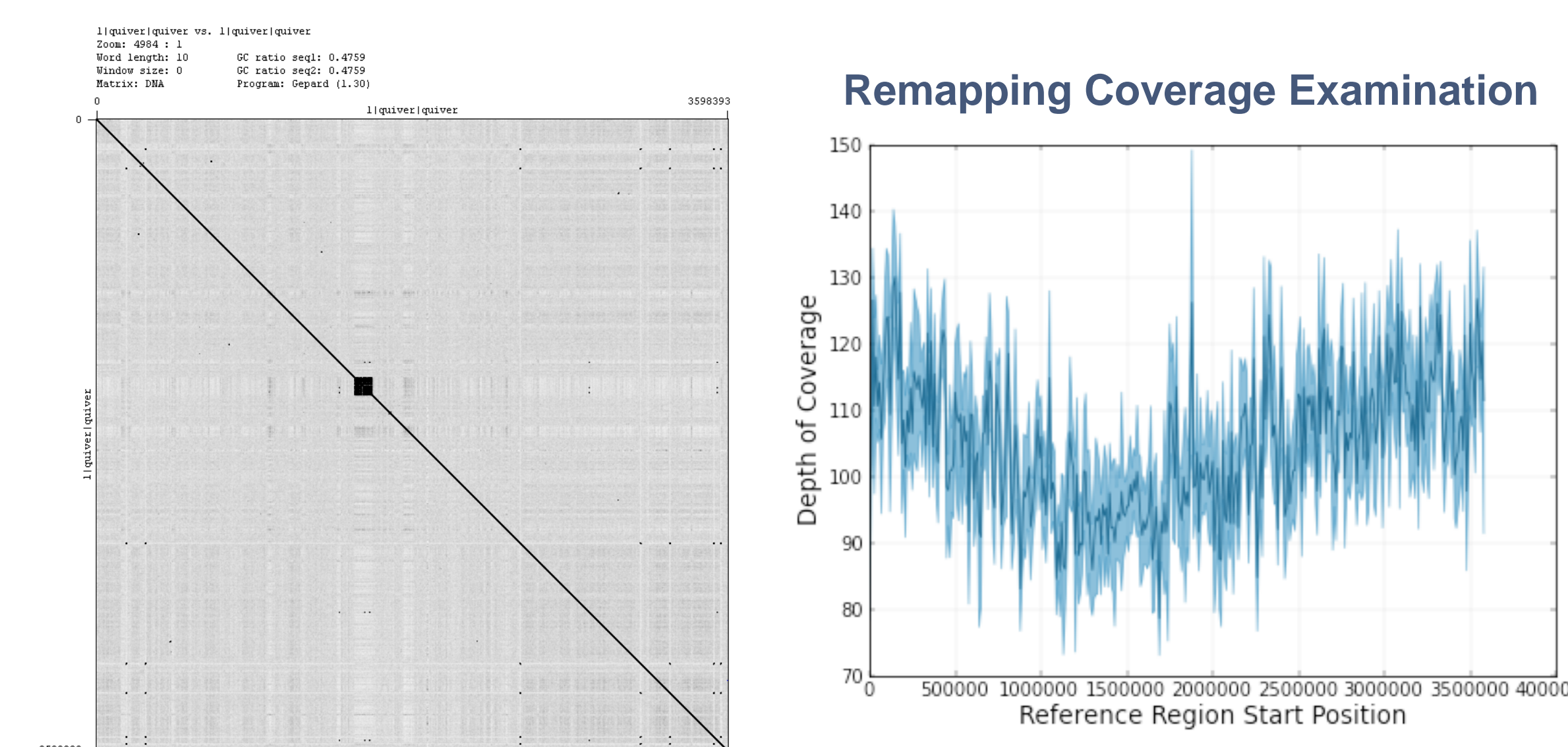


DATA ANALYSIS AND RESULTS

Data analysis with HGAP (de novo assembly using TdT read data only) and subsequent SMRT analysis for base methylation detection.

Two examples are presented here:

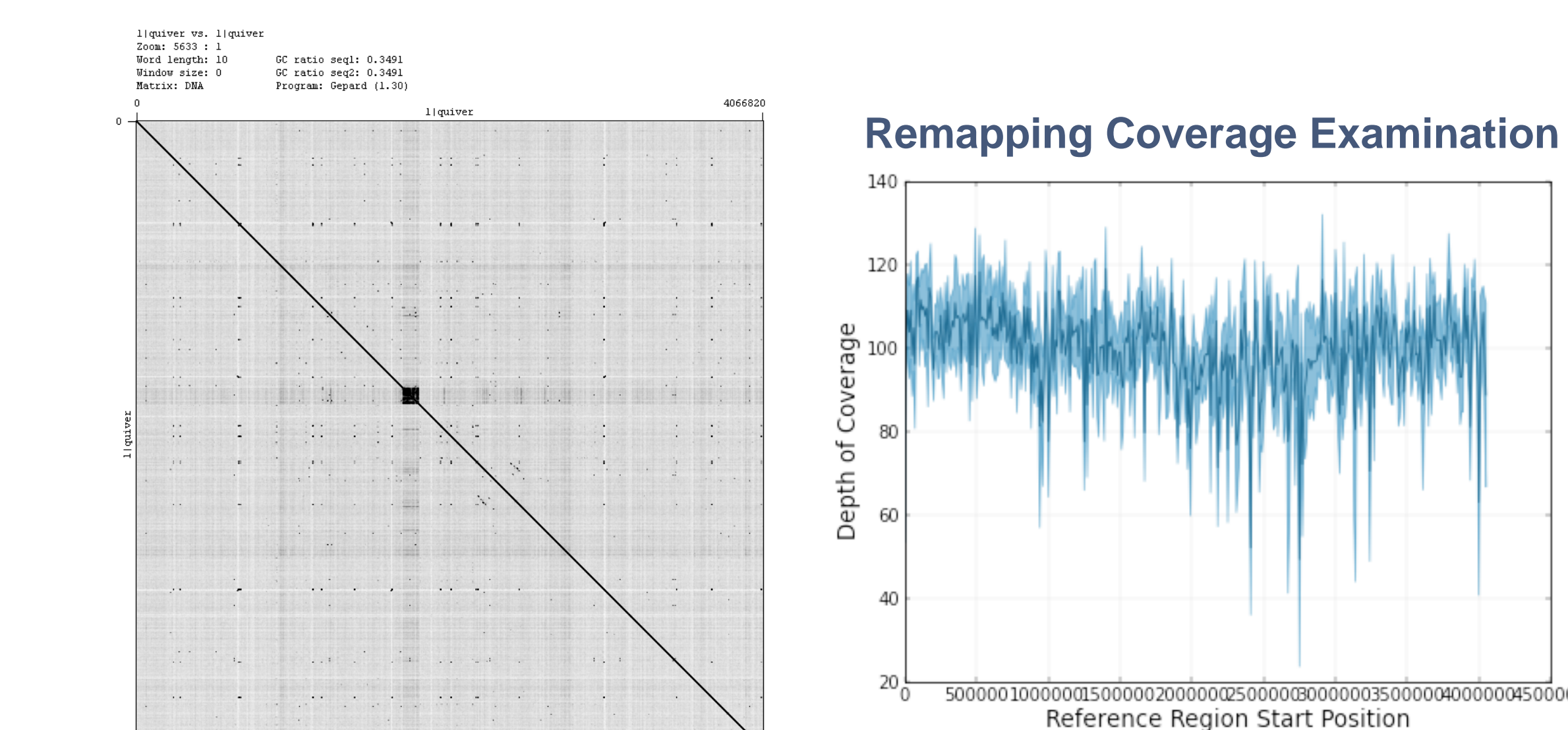
Tolomonas sp. BRL6-1: HGAP produced 1 contig with 3,598,394 bases.



Remapping Quality Examination

	# of Post-Filter Reads	# of Mapped Reads	# of Mapped Bases	Mean Mapped Read Length	95th Percentile Mapped Read Length	Maximum Mapped Read Length	Mean Mapped Subread Accuracy	# of Mapped Subreads	Mean Mapped Subread Length	Mean Mapped Full Subread Length	Mean Max Subread Length
All Movies	163612	146373	39862957 bp	2724 bp	6826 bp	15893 bp	84.75%	146507	2721 bp	0 bp	2722 bp

Gillisia sp. JM1: HGAP produced 1 contig with 4,066,858 bases.

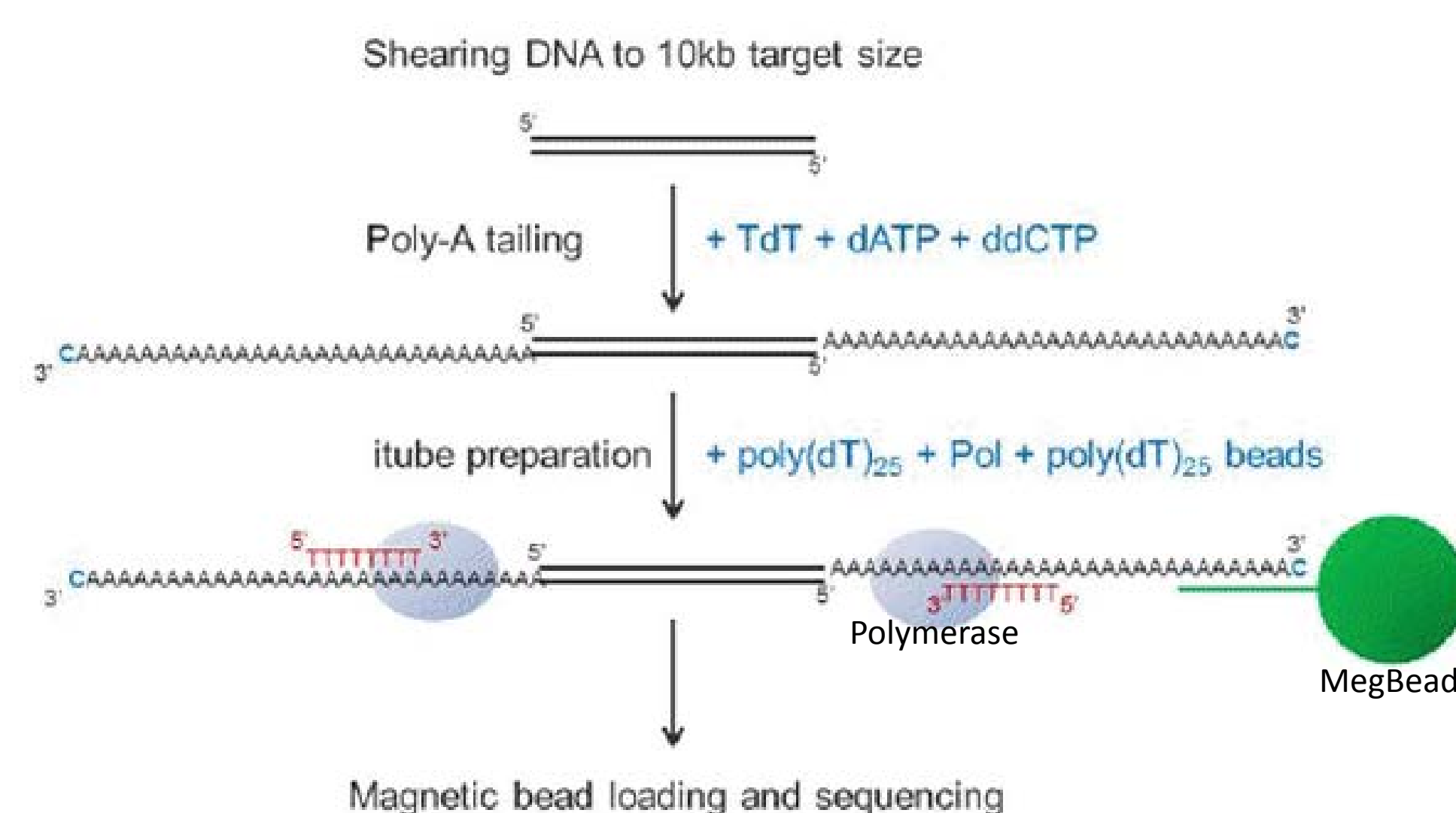


Remapping Quality Examination

	# of Post-Filter Reads	# of Mapped Reads	# of Mapped Bases	Mean Mapped Read Length	95th Percentile Mapped Read Length	Maximum Mapped Read Length	Mean Mapped Subread Accuracy	# of Mapped Subreads	Mean Mapped Subread Length	Mean Mapped Full Subread Length	Mean Max Subread Length
All Movies	187417	167027	430940225 bp	2580 bp	6365 bp	17104 bp	84.21%	167310	2575 bp	332 bp	2578 bp

m6A methylation motifs were also detected in both genomes.

PRINCIPLE OF THE METHOD



Standard SMRTbell 10kb lib	TdT 10kb lib
5-10ug	100-200ng
Shear DNA to 10kb (Qubit & Bioanalyzer) 1.5h	Shear DNA to 10kb (Qubit) 1h
Damage & end Repair 1h	Damage Repair 1h
Blunt end ligation & Exo treatment (Qubit & Bioanalyzer) 3.5h	Exo treatment & Poly-A tailing (Qubit) 2h
Total Prep Time 6h	Total Prep Time 4h
Lib Yield 15-25%	Lib Yield 70-80%