



Sujin Kim¹, Seung-Chul Shin², Hyun Park² & Jong-Eun Lee¹ ¹DNA Link, Inc. Seoul Korea / ²Korea Polar Research Institute, Incheon Korea

Abstract

In the last year, high-throughput sequencing technologies have progressed from proof-of-concept to production quality. Although each technology is able to produce vast quantities of sequence information, in every case the underlying chemistry limits reads to very short lengths. We present a examining de novo assembly comparison with bacterial genome assembly varying genome size (from 3.1Mb to 7.6Mb) and different G+C contents (from 43% to 71%), respectively. We analyzed Solexa reads, 454 reads and Pacbio RS reads from Streptomyces sp. (Genome size, 7.6 Mb; G+C content, 71%), Psychrobacter sp. (Genome size, 3.5 Mb; G+C content, 43%), Salinibacterium sp. (Genome size, 3.1 Mb; G+C content, 61%) and Frigoribacterium sp. (Genome size, 3.3 Mb; G+C content, 63%). We assembly each bacterial genome using Celera assembler 7.0 with and without PacBio RS reads. We found out that the assemble result with Pacbio RS reads have less contigs and scaffolds, and better N50 values.

Summary

General Infomation

Strain No.	Scientific name	GC content (%)	Genome Size (Mb)	Hiseq (Coverage)	PBcR (Coverage)	GSFLX (Coverage)	CCS (Coverage)
PAMC 26508	Streptomyces sp.	70.09	7.6	100X	15X	6X	27X
PAMC 26555	Frigoribacterium sp.	62.87	3.3	100X	26X	20X	
PAMC 21119	Psychrobacter sp.	43.34	3.5	60X	14X	3X	
PAMC 21357	Salinibacterium sp.	60.2	3.1	100X	16X	8X	

Assembly Summary

Strain No. Scientific name		PAMC 26508	PAMC 26555	PAMC 21119	PAMC 21357
		Streptomyces sp.	Frigoribacterium sp.	Psychrobacter sp.	Salinibacterium sp.
	Hiseq	185	93	76	26
Contigs (EA)	PBcR	26	60	39	7
	PBcR+GSFLX	36	96	67	9
Max Contig Bases Bases (bp)	Hiseq	221,220	366,603	255,324	643,647
	PBcR	1,915,364	1,164,979	582,027	2,110,572
	PBcR+GSFLX	1,869,455	1,492,435	598,064	2,112,607
	Hiseq	68,326	136,524	82,883	233,362
N50 Contig Bases Bases (bp)	PBcR	1,268,506	263,283	511,916	2,110,572
	PBcR+GSFLX	1,434,415	469,317	258,230	2,112,607
	Hiseq	141	38	60	20
Total Big Contigs (10kb <)	PBcR	11	16	16	3
	PBcR+GSFLX	11	9	19	3
	Hiseq	7,372,844	3,203,637	3,372,293	3,082,467
Big Contig Length	PBcR	7,637,132	3,248,689	3,439,229	3,111,472
	PBcR+GSFLX	7,723,210	3,246,808	3,441,479	3,110,513





Genome sequencing of microbial genomes using Single Molecule Real-time sequencing (SMRT) technology



	Illumina (60X)	PBcR (14X)	PBcR(16X)+GSFLX(3X)
Scaffolds			
TotalScaffolds (EA)	67	39	60
TotalContigsInScaffolds (EA)	76	39	67
TotalBasesInScaffolds (bp)	3,443,158	3,502,463	3,591,719
MaxBasesInScaffolds (bp)	352,709	582,027	598,064
N50ScaffoldBases (bp)	93,139	511,916	258,230

Contigs		26	
TotalContigsInScaffolds (EA)	76	39	67
MaxContigLength (bp)	255,324	582,027	598,064
N50ContigBases (bp)	82,883	511,916	258,230
TotalBigContigs (10kb<)	60	16	19
BigContigLength (bp)	3,372,293	3,439,229	3,441,479

PAMC21357 : Salinibacterium sp. GC content : 60.2 %

Genome Size : 3.1 Mb

Assembler : Celera Assembler Sequening technologies Illumina : 500bp paired-end library GS-FLX : 8kb paired-end library Pacbio : Continuous Long Read

Assembly results

	Illumina (100X)	PBcR (16X)	PBcR(16X)+GSFLX(8X)
Scaffolds			
TotalScaffolds (EA)	9	7	9
TotalContigsInScaffolds (EA)	26	7	9
TotalBasesInScaffolds (bp)	3,102,350	3,119,401	3,120,661
MaxBasesInScaffolds (bp)	761,436	2,110,572	2,112,607
N50ScaffoldBases (bp)	631,463	2,110,572	2.112.607

Contigs			
TotalContigsInScaffolds (EA)	26	7	9
MaxContigLength (bp)	643,647	2,110,572	2,112,607
N50ContigBases (bp)	233,362	2,110,572	2,112,607
TotalBigContigs (10kb<)	20	3	3
BigContigLength (bp)	3,082,467	3,111,472	3,110,513

The improvement of de novo assemblies with Pacbio RS



Assembly comparison





