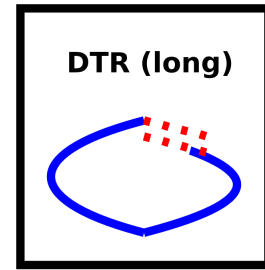
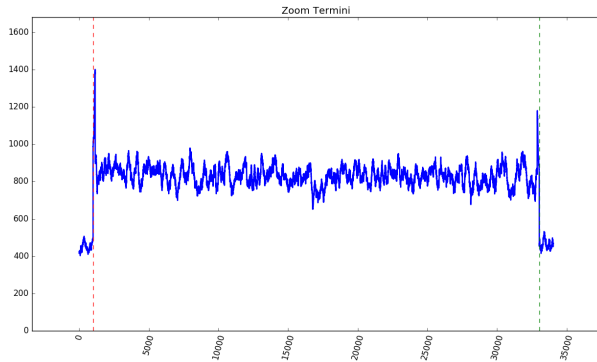


# XaC1 PhageTerm Analysis



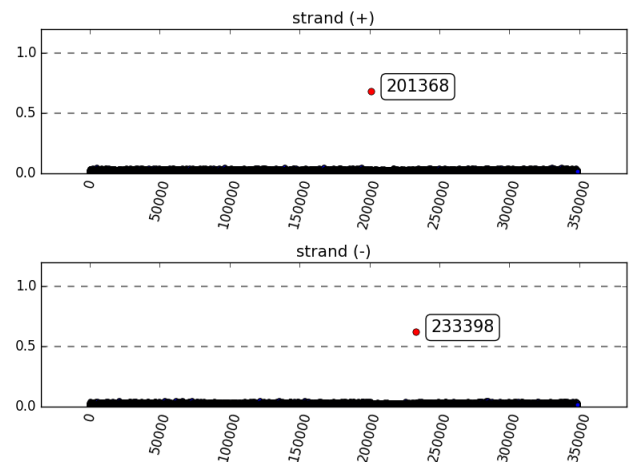
## PhageTerm Method

Ends	Left (red)	Right (green)	Permuted	Orientation	Class	Type
Redundant	201368	233398	No	NA	DTR (long)	T5

\*Direct Terminal Repeats: 32031 bp

Strand	Location	T	pvalue
+	201368	0.68	2.93e-286
	5130	0.05	1.00e+00
	96913	0.05	1.00e+00
	139661	0.05	1.00e+00
	194553	0.05	1.00e+00
-	233398	0.62	4.93e-167
	283605	0.05	1.00e+00
	21281	0.05	1.00e+00
	122038	0.05	1.00e+00
	284918	0.05	1.00e+00

## T (Start. Pos. Cov. / Whole Cov.)



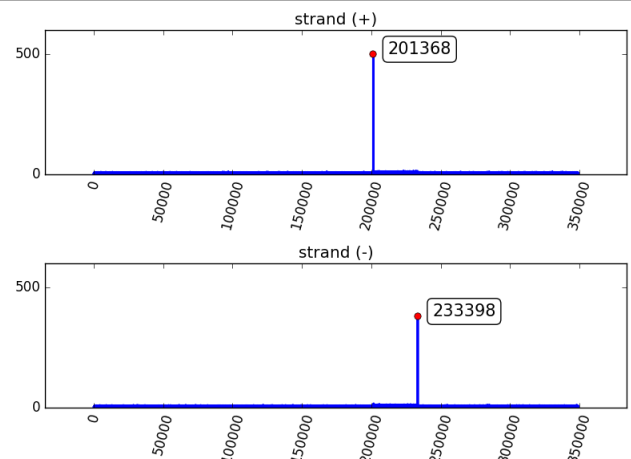
## Li's Method

Packaging	Termini	Forward	Reverse	Orientation
COS	Fixed	Obvious Termini	Obvious Termini	Forward

\*Direct Terminal Repeats: 32031 bp

Strand	Location	SPC	R
+	201368	502	36.0
	230298	14	-
	223196	14	-
	233171	13	-
	232104	13	-
-	233398	381	24.0
	201618	16	-
	205722	14	-
	201588	14	-
	228242	13	-

## SPC

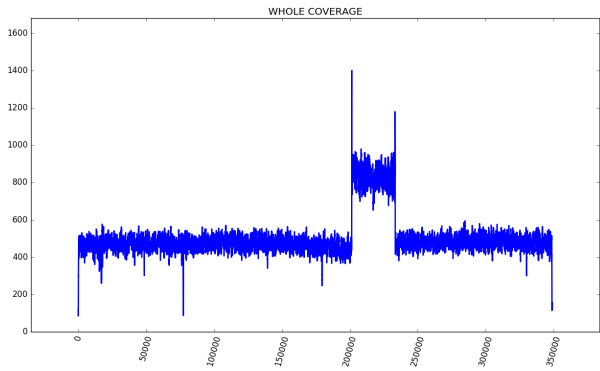


Analysis Methodology

PhageTerm software uses raw reads of a phage sequenced with a sequencing technology using random fragmentation and its genomic reference sequence to determine the termini position. The process starts with the alignment of NGS reads to the phage genome in order to calculate the starting position coverage (SPC), where a hit is given only to the position of the first base in a successfully aligned read (the alignment algorithm uses the lenght of the seed (default: 20) for mapping and does not accept gap or mismatch to speed up the process). Then the program apply 2 distinct scoring methods: i) a statistical approach based on the Gamma law; and ii) a method derived from LI and al. 2014 paper.

General set-up and mapping informations

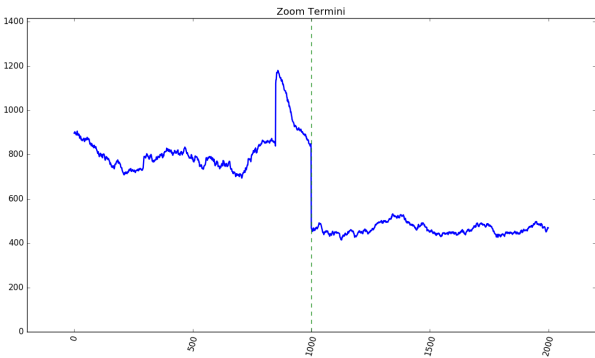
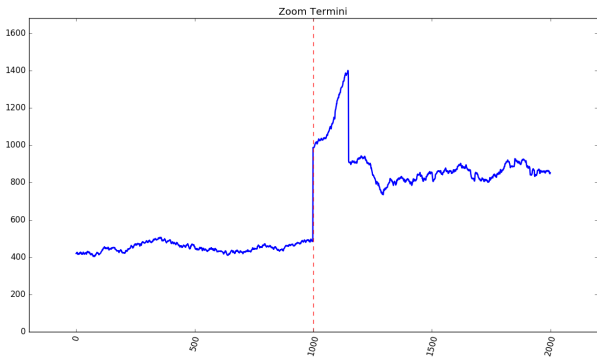
Phage Genome	348967 bp
Sequencing Reads	19839651
Mapping Reads	97 %
OPTIONS	
Mapping Seed	20
Surrounding	20
Host Analysis	No



Highest peak of each side coverage graphics

Whole Coverage Zoom (Left)

Whole Coverage Zoom (Right)



General controls information

Whole genome coverage	252	OK
Weak genome coverage	0.0 %	OK
Reads lost during alignment	2.7 %	OK

i) PhageTerm method

Reads are mapped on the reference to determine the starting position coverage (SPC) as well as the coverage (COV) in each orientation. These values are then used to compute the variable  $T = SPC / COV$ . The average value of  $T$  at positions along the genome that are not termini is expected to be  $1/F$ , where  $F$  is the average fragment size. For the termini that depends of the packaging mode. Cos Phages: no reads should start before the terminus and therefore  $X=1$ . DTR phages: for  $N$  phages present in the sample, there should be  $N$  fragments that start at the terminus and  $N$  fragments that cover the edge of the repeat on the other side of the genome as a results  $T$  is expected to be 0.5. Pac phages: for  $N$  phages in the sample, there should be  $N/C$  fragments starting at the pac site, where  $C$  is the number of phage genome copies per concatemer. In the same sample  $N$  fragments should cover the pac site position,  $T$  is expected to be  $(N/C)/(N+N/C) = 1/(1+C)$ . To assess whether the number of reads starting at a given position along the genome can be considered a significant outlier, PhageTerm first segments the genome according to coverage using a regression tree. A gamma distribution is then fitted to SPC for each segment and an adjusted  $p$ -value is computed for each position. Finally if several significant peaks are detected within a small sequence window (default: 20bp), their  $T$  values are merged.

Nearby Termini (Forward / Reverse)	0 / 0	Peaks localized 20 bases around the maximum
------------------------------------	-------	---

ii) Li's method

The second approach is based on the calculation and interpretation of three specific ratios  $R1$ ,  $R2$  and  $R3$  as suggested in a previous publication from Li et al. 2014. The first ratio, is calculated as follow: the highest starting frequency found on either the forward or reverse strands is divided by the average starting frequency,  $R1 = (\text{highest frequency} / \text{average frequency})$ . Li's et al. have proposed three possible interpretation of the  $R1$  ratio. First, if  $R1 < 30$ , the phage genome does not have any termini, and is either circular or completely permuted and terminally redundant. The second interpretation for  $R1$  is when  $30 \leq R1 \leq 100$ , suggesting the presence of preferred termini with terminal redundancy and apparition of partially circular permutations. At last if  $R1 > 100$  that is an indication that at least one fixed termini is present with terminase recognizing a specific site. The two other ratios are  $R2$  and  $R3$  and the calculation is done in a similar manner.  $R2$  is calculated using the highest two frequencies ( $T1-F$  and  $T2-F$ ) found on the forward strand and  $R3$  is calculated using the highest two frequencies ( $T1-R$  and  $T2-R$ ) found on the reverse strand. To calculate these two ratios, we divide the highest frequency by the second highest frequency  $T2$ . So  $R2 = (T1-F / T2-F)$  and  $R3 = (T1-R / T2-R)$ . These two ratios are used to analyze termini characteristics on each strand taken individually. Li et al. suggested two possible interpretations for  $R2$  and  $R3$  ratios combine to  $R1$ . When  $R1 < 30$  and  $R2 < 3$ , we either have no obvious termini on the forward strand, or we have multiple preferred termini on the forward strand, if  $30 \leq R1 \leq 100$ . If  $R2 > 3$ , it is suggested that there is an obvious unique termini on the forward strand. The same reasoning is applicable for the result of  $R3$ . Combining the results for ratios found with this approach, it is possible to make the first prediction for the viral packaging mode of the analyzed phage. A unique obvious termini present at both ends (both  $R2$  and  $R3 > 3$ ) reveals the presence of a COS mode of packaging. The headful mode of packaging PAC is concluded when we have a single obvious termini only on one strand. A whole coverage around 500X is needed for this method to be reliable.

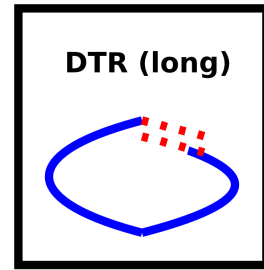
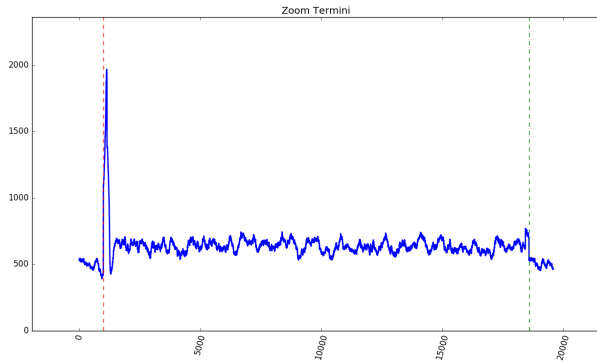
Nearby Termini (Forward / Reverse)	0 / 1	Peaks localized 20 bases around the maximum
R1 - highest freq./average freq.	287	At least one fixed termini is present with terminase recognizing a specific site.
R2 Forw - highest freq./second freq.	36	Unique termini on the forward strand.
R3 Rev - highest freq./second freq.	24	Unique termini on the reverse strand.

Please cite: Sci. Rep. DOI 10.1038/s41598-017-07910-5

Garneau, Depardieu, Fortier, Bikard and Monot. PhageTerm: Determining Bacteriophage Termini and Packaging using NGS data.

Report generated : Wed Mar 3 05:55:24 2021

# XbC2 PhageTerm Analysis



## PhageTerm Method

Ends	Left (red)	Right (green)	Permuted	Orientation	Class	Type
Redundant	317931	335511	No	NA	DTR (long)	T5

\*Direct Terminal Repeats: 17581 bp

Strand	Location	T	pvalue	T (Start. Pos. Cov. / Whole Cov.)
+	317931	0.72	5.56e-168	
	260885	0.05	1.00e+00	
	51260	0.05	1.00e+00	
	365010	0.05	1.00e+00	
	216506	0.04	1.00e+00	
-	335511	0.39	1.04e-43	
	318195	0.14	2.52e-13	
	115519	0.05	1.00e+00	
	318187	0.05	4.67e-03	
	318192	0.05	8.47e-02	

## Li's Method

Packaging	Termini	Forward	Reverse	Orientation
COS	Fixed	Obvious Termini	Obvious Termini	Forward

\*Direct Terminal Repeats: 17581 bp

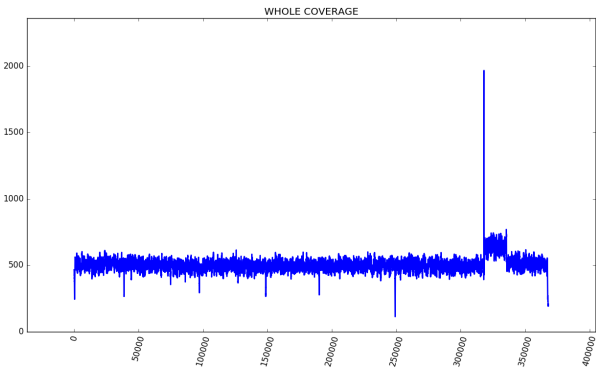
Strand	Location	SPC	R	SPC
+	317931	646	54.0	
	365010	12	-	
	333610	12	-	
	261973	12	-	
	51260	12	-	
-	335511	152	3.0	
	318195	59	-	
	318185	28	-	
	318187	27	-	
	318181	24	-	

Analysis Methodology

PhageTerm software uses raw reads of a phage sequenced with a sequencing technology using random fragmentation and its genomic reference sequence to determine the termini position. The process starts with the alignment of NGS reads to the phage genome in order to calculate the starting position coverage (SPC), where a hit is given only to the position of the first base in a successfully aligned read (the alignment algorithm uses the lenght of the seed (default: 20) for mapping and does not accept gap or mismatch to speed up the process). Then the program apply 2 distinct scoring methods: i) a statistical approach based on the Gamma law; and ii) a method derived from Li and al. 2014 paper.

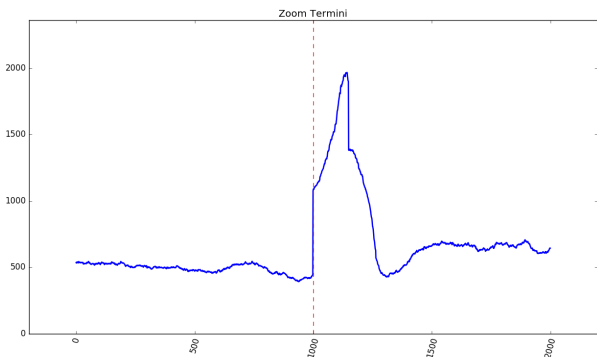
General set-up and mapping informations

Phage Genome	367901 bp
Sequencing Reads	17043954
Mapping Reads	97 %
OPTIONS	
Mapping Seed	20
Surrounding	20
Host Analysis	No

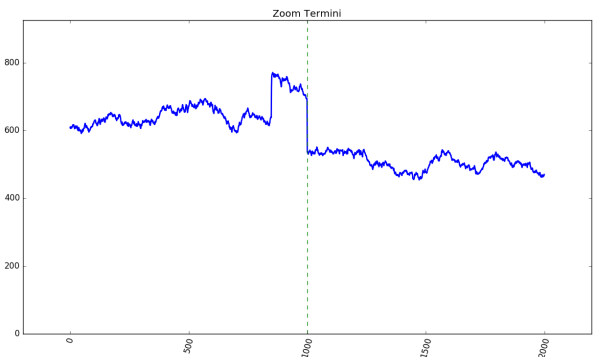


Highest peak of each side coverage graphics

Whole Coverage Zoom (Left)



Whole Coverage Zoom (Right)



General controls information

Whole genome coverage	252	OK
Weak genome coverage	0.0 %	OK
Insert mean size	488	Mean insert estimated from paired-end reads
Reads lost during alignment	3.0 %	OK

i) PhageTerm method

Reads are mapped on the reference to determine the starting position coverage (SPC) as well as the coverage (COV) in each orientation. These values are then used to compute the variable  $T = SPC / COV$ . The average value of  $T$  at positions along the genome that are not termini is expected to be  $1/F$ , where  $F$  is the average fragment size. For the termini that depends of the packaging mode. Cos Phages: no reads should start before the terminus and therefore  $X=1$ . DTR phages: for  $N$  phages present in the sample, there should be  $N$  fragments that start at the terminus and  $N$  fragments that cover the edge of the repeat on the other side of the genome as a results  $T$  is expected to be 0.5. Pac phages: for  $N$  phages in the sample, there should be  $N/C$  fragments starting at the pac site, where  $C$  is the number of phage genome copies per concatemer. In the same sample  $N$  fragments should cover the pac site position,  $T$  is expected to be  $(N/C)/(N+N/C) = 1/(1+C)$ . To assess whether the number of reads starting at a given position along the genome can be considered a significant outlier, PhageTerm first segments the genome according to coverage using a regression tree. A gamma distribution is then fitted to SPC for each segment and an adjusted  $p$ -value is computed for each position. Finally if several significant peaks are detected within a small sequence window (default: 20bp), their  $T$  values are merged.

Nearby Termini (Forward / Reverse)	0 / 0	Peaks localized 20 bases around the maximum
------------------------------------	-------	---

ii) Li's method

The second approach is based on the calculation and interpretation of three specific ratios R1, R2 and R3 as suggested in a previous publication from Li et al. 2014. The first ratio, is calculated as follow: the highest starting frequency found on either the forward or reverse strands is divided by the average starting frequency,  $R1 = (\text{highest frequency} / \text{average frequency})$ . Li's et al. have proposed three possible interpretation of the R1 ratio. First, if  $R1 < 30$ , the phage genome does not have any termini, and is either circular or completely permuted and terminally redundant. The second interpretation for R1 is when  $30 \leq R1 \leq 100$ , suggesting the presence of preferred termini with terminal redundancy and apparition of partially circular permutations. At last if  $R1 > 100$  that is an indication that at least one fixed termini is present with terminase recognizing a specific site. The two other ratios are R2 and R3 and the calculation is done in a similar manner. R2 is calculated using the highest two frequencies (T1-F and T2-F) found on the forward strand and R3 is calculated using the highest two frequencies (T1-R and T2-R) found on the reverse strand. To calculate these two ratios, we divide the highest frequency by the second highest frequency T2. So  $R2 = (T1-F / T2-F)$  and  $R3 = (T1-R / T2-R)$ . These two ratios are used to analyze termini characteristics on each strand taken individually. Li et al. suggested two possible interpretations for R2 and R3 ratios combine to R1. When  $R1 < 30$  and  $R2 < 3$ , we either have no obvious termini on the forward strand, or we have multiple preferred termini on the forward strand, if  $30 \leq R1 \leq 100$ . If  $R2 > 3$ , it is suggested that there is an obvious unique termini on the forward strand. The same reasoning is applicable for the result of R3. Combining the results for ratios found with this approach, it is possible to make the first prediction for the viral packaging mode of the analyzed phage. A unique obvious termini present at both ends (both R2 and R3 > 3) reveals the presence of a COS mode of packaging. The headful mode of packaging PAC is concluded when we have a single obvious termini only on one strand. A whole coverage around 500X is needed for this method to be reliable.

Nearby Termini (Forward / Reverse)	1 / 0	Peaks localized 20 bases around the maximum
R1 - highest freq./average freq.	365	At least one fixed termini is present with terminase recognizing a specific site.
R2 Forw - highest freq./second freq.	54	Unique termini on the forward strand.
R3 Rev - highest freq./second freq.	3	Unique termini on the reverse strand.

Please cite: Sci. Rep. DOI 10.1038/s41598-017-07910-5

Garneau, Depardieu, Fortier, Bikard and Monot. PhageTerm: Determining Bacteriophage Termini and Packaging using NGS data.

Report generated : Tue Feb 23 20:43:35 2021

**Name:** XaC1**Conformation:** linear**Enzymes:** EcoRI, EcoRV, NdeI**Noncutters:**

Name	Sequence	Site Length	Overhang	Frequency	Cut Positions
<a href="#">EcoRV</a>	GATATC	6	blunt	27	5062, 5080, 18506, 18800, 19760, 20228, 25195, 26593, 30188, 46366, 70650, 106976, 164169, 167306, 169111, 214413, 227544, 237933, 248954, 276996, 283085, 307177, 314479, 332609, 335230, 346892, 348798
<a href="#">EcoRI</a>	GAATTC	6	five_prime	140	34650, 35361, 35907, 37586, 41330, 45461, 45763, 46458, 49914, 51452, 52621, 53443, 53473, 53919, 61826, 62888, 67950, 69226, 71310, 73384, 74711, 74980, 75867, 76218, 77678, 79156, 80241, 81299, 81527, 82978, 83323, 92478, 106734, 107963, 108319, 108757, 114943, 117798, 117963, 118286, 119377, 119913, 122657, 123361, 125097, 125409, 125985, 127832, 129180, 133198, 134083, 134280, 135980, 138617, 149187, 149382, 151780, 161858, 166744, 176871, 178190, 179280, 179739, 180221, 187376, 187512, 188463, 189686, 190332, 190535, 191431, 193118, 193292, 195741, 197132, 203303, 208057, 208396, 209653, 209744, 211217, 217967, 225300, 226152, 230601, 232763, 235107, 240908, 241063, 243694, 244587, 244599, 244704, 246968, 250398, 250982, 255155, 255530, 255555, 255712, 256638, 261549, 262998, 264765, 269447, 269669, 270917, 272307, 274225, 279081, 281757, 284810, 285802, 288349, 289070, 289363, 293516, 295162, 303541, 303751, 303861, 304434, 307121, 307193, 311303, 318669, 320775, 322297, 325841, 326766, 326951, 327879, 330406, 331598, 332007, 333392, 336345, 338801, 339377, 340651
<a href="#">NdeI</a>	CATATG	6	five_prime	162	1679, 6786, 18572, 19138, 21660, 24177, 27176, 28754, 28775, 30673, 30955, 33258, 33896, 38922, 39162, 42597, 42784, 44821, 48362, 49628, 50231, 52551, 53609, 55484, 56435, 59611, 60150, 61094, 61796, 63342, 64394, 69838, 70788, 71327, 72657, 72861, 73416, 80664, 81898, 85886, 87412, 89190, 89982, 93084, 95386, 97130, 99904, 100379, 103380, 106463, 106818, 109597, 110277, 110940, 111650, 114391, 114913, 115015, 115384, 118219, 118849, 119020, 120368, 123231, 125336, 128211, 128903, 129239, 130079, 132177, 136015, 137661, 139513, 139630, 146741, 153045, 153800, 156388, 163361, 163687, 164020, 168336, 168710, 171412, 178385, 179933, 180784, 182128, 191333, 192266, 195602, 195638, 195647, 198646, 199384, 201341, 202043, 203845, 208104, 208797, 208920, 209827, 211385, 213917, 214821, 225434, 228064, 229017,

Name	Sequence	Site Length	Overhang	Frequency	Cut Positions
					238536, 238707, 241887, 241985, 242126, 247002, 247094, 251146, 252235, 261800, 262541, 264073, 268314, 268386, 271832, 271922, 272600, 274286, 277446, 277812, 279945, 281749, 282272, 282437, 283520, 286009, 287685, 287712, 288519, 291680, 295198, 296007, 296543, 300117, 300837, 301821, 302821, 303623, 303931, 312777, 312837, 314107, 315087, 318423, 319511, 320067, 324287, 335065, 336443, 337480, 337570, 337810, 338289, 342982



**Name:** XbC2**Conformation:** linear**Enzymes:** EcoRI, EcoRV, NdeI**Noncutters:**

Name	Sequence	Site Length	Overhang	Frequency	Cut Positions
<a href="#">EcoRV</a>	GATATC	6	blunt	9	4388, 8809, 83399, 139105, 179294, 258149, 261739, 342285, 367722
<a href="#">NdeI</a>	CATATG	6	five_prime	148	3254, 12730, 13565, 13964, 14970, 17645, 20533, 22941, 24718, 25727, 28474, 31002, 33568, 34770, 35327, 42158, 43533, 53416, 54324, 56051, 56441, 58205, 59459, 61268, 63841, 69452, 69537, 72279, 76958, 78086, 79165, 79172, 80615, 80771, 81242, 82148, 83207, 83555, 85282, 87487, 88077, 100475, 106637, 110534, 112497, 112706, 113791, 115874, 119191, 119332, 132048, 134852, 135341, 141510, 144998, 150537, 153936, 154191, 154326, 155493, 155907, 156492, 157339, 158136, 158250, 159474, 161617, 169033, 169090, 174262, 174337, 174429, 179160, 182293, 185592, 188417, 189526, 190758, 192650, 193078, 197057, 198284, 198984, 202475, 202577, 203169, 206900, 209144, 230225, 236119, 236651, 238376, 239891, 241128, 241317, 241674, 242073, 242424, 247634, 248509, 257750, 259419, 259551, 261649, 262229, 265074, 265330, 266163, 267251, 268891, 269110, 275402, 277939, 279767, 280447, 282639, 289375, 292997, 298069, 300973, 305815, 317350, 319828, 321152, 322755, 325071, 326310, 329993, 334993, 335722, 336232, 338053, 339981, 340089, 344430, 347115, 350554, 350979, 352989, 354294, 354660, 357638, 359456, 364960, 366190, 366257, 366704, 366902
<a href="#">EcoRI</a>	GAATTC	6	five_prime	199	8593, 9912, 14527, 14957, 17685, 22107, 22613, 23043, 24064, 24399, 25772, 27381, 27675, 30973, 34883, 37306, 37933, 38993, 39705, 39849, 40388, 41307, 43508, 44933, 45197, 47135, 49388, 51281, 51834, 53366, 53648, 53847, 55390, 57194, 61861, 64618, 64856, 66517, 66553, 67177, 68539, 68802, 70841, 71830, 73334, 74553, 76022, 76596, 79748, 82183, 82488, 83998, 85224, 87240, 89210, 89528, 91473, 95842, 97146, 97870, 99309, 100426, 101292, 103356, 105267, 111953, 112588, 113412, 114050, 116471, 120273, 125126, 125282, 127555, 128266, 128485, 129295, 134063, 134476, 134606, 134765, 138290, 139999, 140933, 142010, 149022, 149034, 149370, 153493, 156522, 158694, 160066, 160291, 165356, 170414, 172079, 179514, 180440, 181615, 181753, 183473, 189266, 190446, 193353, 194371, 194413, 195418, 196944, 202008, 204994, 205776, 206575, 206589, 207439, 210171, 212308, 212811, 215595, 216399, 218954, 225356, 225520,

Name	Sequence	Site Length	Overhang	Frequency	Cut Positions
					225668, 226687, 228796, 228858, 230299, 230817, 232812, 233560, 237288, 237688, 238642, 242690, 243292, 246426, 247752, 247911, 248933, 251480, 253104, 257559, 257724, 258726, 260167, 261432, 262692, 263128, 263485, 265289, 266157, 267099, 269357, 269766, 270652, 270667, 271189, 271226, 273576, 273936, 274647, 275107, 278964, 283434, 284256, 286234, 286937, 289790, 292561, 293722, 295062, 295572, 296922, 300370, 301979, 308529, 309058, 309184, 312165, 314431, 316667, 316712, 324736, 328566, 328605, 331451, 336686, 337676, 339965, 340753, 342346, 343993, 344128, 345226, 348212, 352197, 353529, 356246, 363871