

# FASTQ QC Report

Report Date	12-21-16
Run ID	161219_D00796_0155_ACAC53ANXX
Project ID	EC-EL-4039
Sample	Sample_OD12_R1
FASTX-Toolkit Version	0.0.13.2
FastQC Version	0.10.1
Dupest Version	0.1.0

This report was automatically generated by the WCMC Epigenomics Core QC pipeline and contains information for assessing the quality of FASTQ sequencing data.

The QC Pipeline executes the following analysis:

1. All FASTQ files for the sample are concatenated to a single file. For paired-end sequences, FASTQ files for each read are concatenated and processed separately, with an "R1" or "R2" appended to the sample name.
2. To identify genomic sequencing bias or low sequence diversity k-length oligonucleotide enrichment is calculated and plotted from the combined FASTQ file using FastQC. *Note:* FastQC only analyses the top 2% of the reads in the FASTQ file and the results are extrapolated over the remainder.
3. Duplication level is estimated from the combined FASTQ file as  $(N - U)/N$  where  $N$  is total reads and  $U$  is the number of unique sequences.
4. Sequencing base call quality statistics are calculated from the combined FASTQ file using FASTX-Toolkit FASTQ Quality Filter.

The report contains the following figures:

1. Sequence Duplication - Estimate of duplication level as a percentage of total reads.
2. Base sequence quality - Calculated from FASTX-Toolkit FASTQ Quality Filter.  
Distribution of base quality scores (Q scores) per sequencing cycle. In a reasonably good sequencing run the majority of the signal should be above Q30. Quality scores are divided into three ranges: green indicates calls of very good quality; orange indicates calls of reasonable quality and red indicates calls of poor quality. Yellow boxes represent the inter-quartile range. Upper and lower whiskers represent the maximum and minimum excluding outliers. The red line represents the median quality and the blue line represents the mean quality.
3. Sequence base content - Percentage of bases represented at each position in the read; calculated from FASTX-Toolkit FASTQ Quality Filter.
4. K-mer content - calculated and plotted by FastQC. From FastQC Help:  
The k-mer analysis checks if there are short fragments of k-length nucleotides that are over represented at certain positions in the reads. In a diversified library there should not be positional bias in its appearance of k-mers. There may be biological reasons why certain k-mers are enriched or depleted overall, but these biases should affect all positions within a sequence equally. In contrast, if certain k-mers are over represented in certain positions then this could indicate issues with library preparation, quality of the input material or sequencing of the adaptors. This analysis measures the number of each 5-mer at each position in the library and then uses a binomial test to look for significant deviations from an even coverage at all positions. Any k-mer with positionally biased enrichment are reported. The top 6 most biased k-mers are additionally plotted to show their distribution. Note that because of the computational overhead associated with calculating k-mer content this analysis is performed on 2% of the reads.
5. Overrepresented sequences - Calculated and plotted by FastQC. From FastQC Help:  
A normal high-throughput library will contain a diverse set of sequences, with no individual sequence making up a tiny fraction of the whole. Finding that a single sequence is very overrepresented in the set either means that it is highly biologically significant, or indicates that the library is contaminated, or not as diverse as you expected.  
This analysis lists all of the reads which make up more than 0.1% of the total. To limit memory use only sequences which appear in the first 200,000 sequences are evaluated for their occurrences in the entire library. It is possible that a sequence which is overrepresented but doesn't appear at the start of the file for some reason

could be missed by this analysis. However, this is unlikely since library preparation and sequencing randomize the genomic elements and therefore the first 200,000 reads are sufficient to represent the diversity in the entire library.

For each overrepresented sequence the program will look for matches in a database of common contaminants and will report the best hit it finds. Hits must be at least 20bp in length and have no more than 1 mismatch. Finding a hit doesn't necessarily mean that this is the source of the contamination, but may provide clues about the true source of contamination. It's also worth pointing out that many adaptor sequences are similar in sequence so a match to an adaptor sequence may not represent the true source of the adaptor.

Because the duplication detection requires an exact sequence match over the whole length of the sequence. Reads over 75bp in length are truncated to 50bp for the purposes of this analysis.

FastQC: <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>

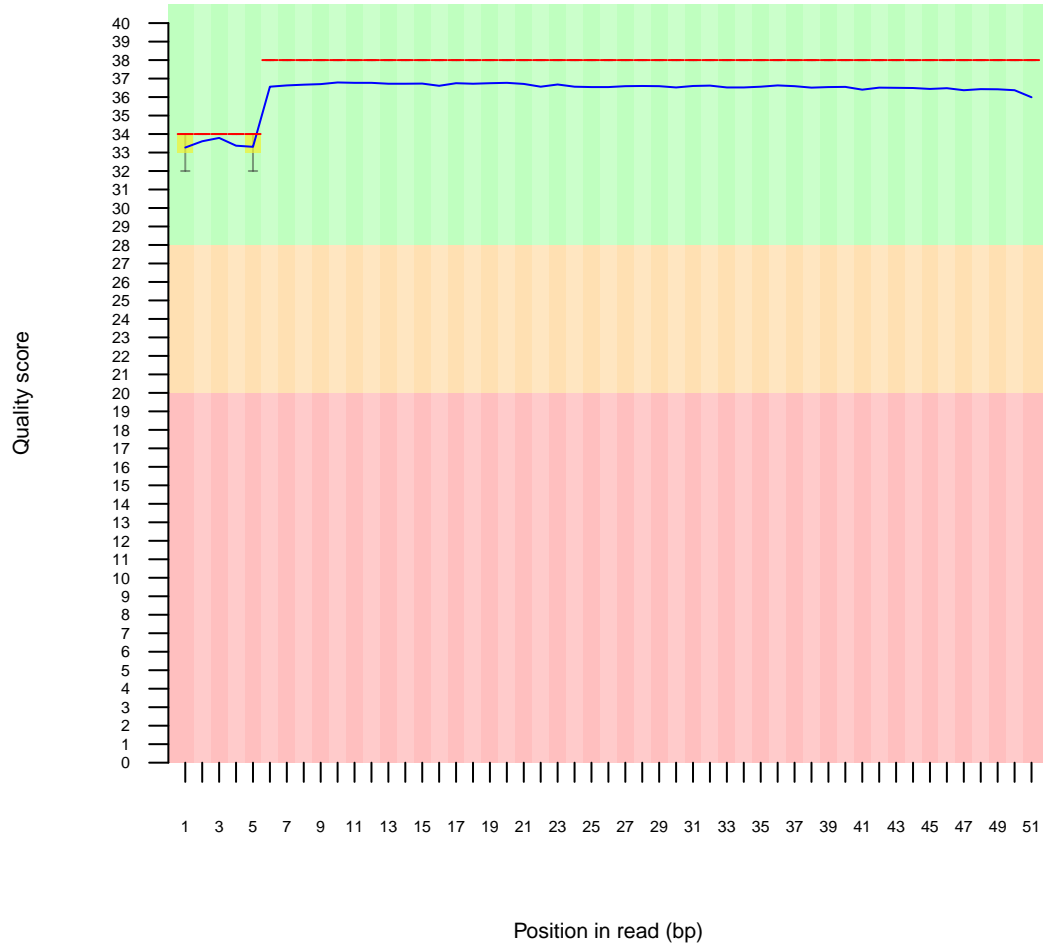
FASTX-Toolkit: [http://hannonlab.cshl.edu/fastx\\_toolkit](http://hannonlab.cshl.edu/fastx_toolkit)

# 1 Sequence Duplication

- Estimated Duplication rate 82.2018%

# 2 Per base sequence quality

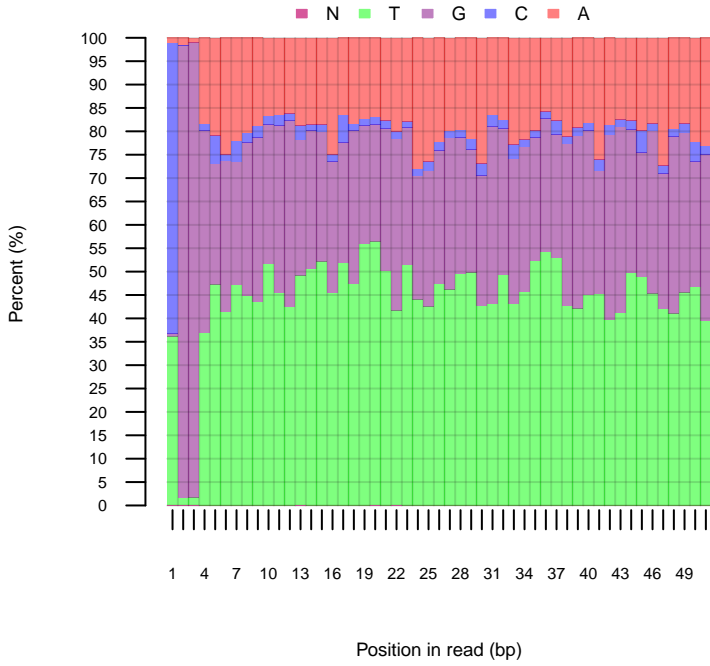
Quality scores across all bases



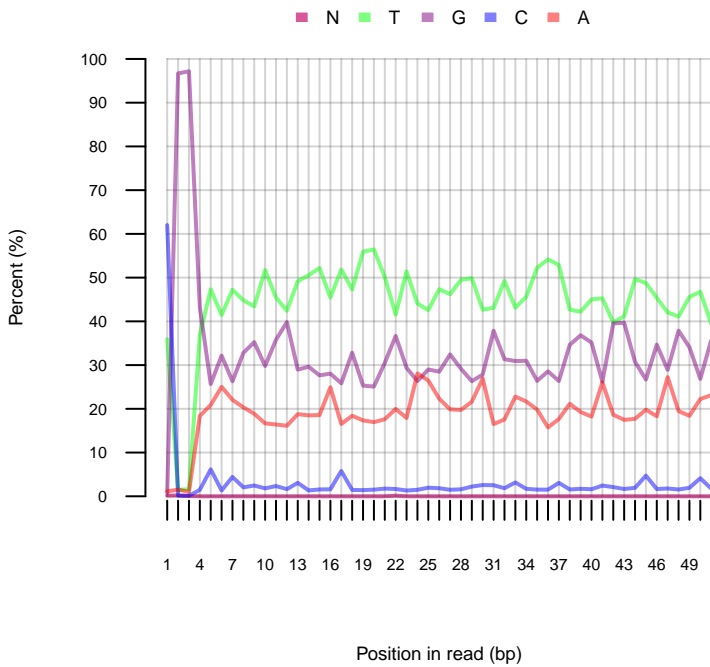
Background colors	Green - calls of very good quality Orange - calls of reasonable quality Red - calls of poor quality
Yellow boxes	Inter-quartile range
Upper and lower whiskers	Maximum and minimum quality excluding outliers
Red line	Median quality
Blue line	Mean quality

### 3 Sequence base content

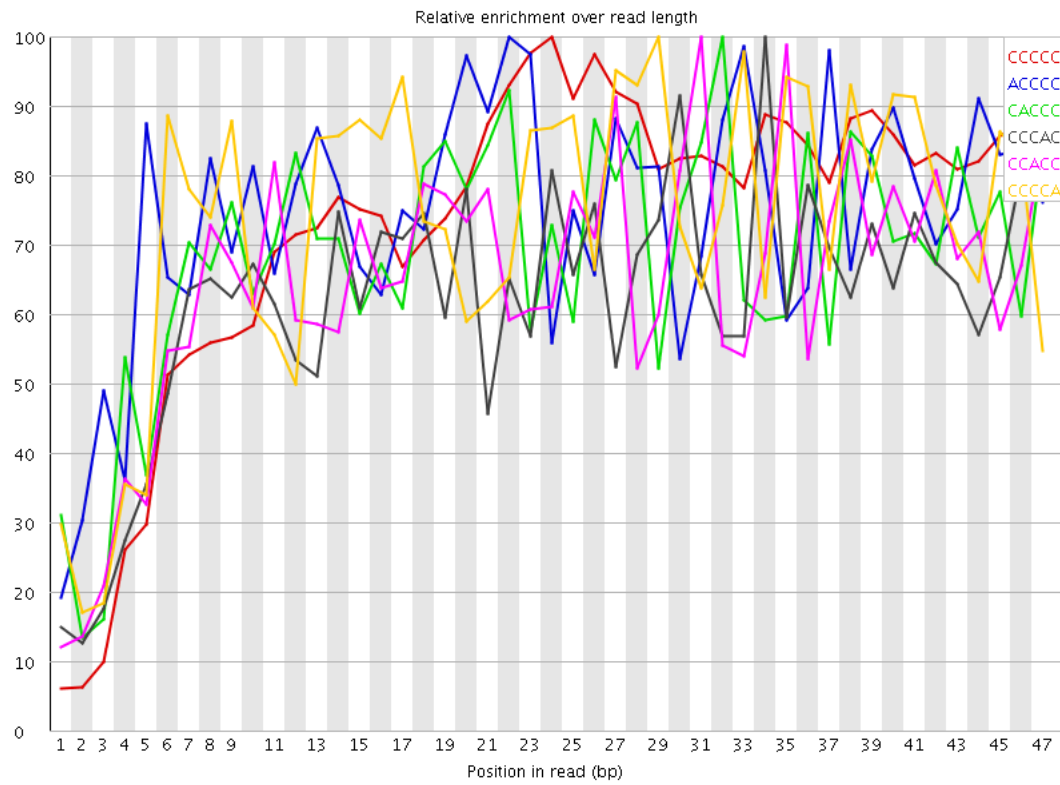
Sequence base content across all positions



Sequence base content across all positions



## 4 Sequence K-mer content



Note: FastQC analyses 2% of the sequence data and results are extrapolated to the rest of the sequence.

Sequence	Count	Obs/Exp Overall	Obs/Exp Max	Max Obs/Exp Position
CCCC	201595	1390.1284	1895.4652	24
ACCC	62775	74.916885	100.510155	22
CACC	62310	74.361946	108.226845	32
CCCAC	60635	72.36297	116.0775	34
CCACC	59950	71.54547	110.189514	31
CCCCA	58360	69.647934	95.048965	29
CACAC	200285	41.367493	273.9721	12
CGGGC	4433295	29.514545	1146.1742	1
CCCGC	28645	19.522497	37.787952	34
GCCCC	28535	19.447529	31.70345	37
CCCGG	27230	18.558132	31.383213	33
CGGCC	27220	18.551313	32.985855	43
CGCCC	26195	17.852743	28.821838	44
CGCGG	2295550	15.282563	435.74234	5
GCGCG	2170105	14.447419	433.63397	4
CTCC	25175	12.8136835	17.69787	28
GGCGC	1781000	11.85696	435.2144	3
CTCCC	22125	11.261281	20.834723	21
CCCTC	22115	11.256191	16.621647	37
TCCCC	21700	11.044963	25.231424	5
CCCTT	20290	10.327294	16.382835	47
CGGCG	1547085	10.299679	294.05536	1
CCACA	39735	8.206992	11.112225	17
CGCGC	120700	8.130277	41.573143	13
CGGAA	3986930	8.044004	245.97914	1
CGGGA	6398630	7.3724914	240.84634	1
ACACC	35655	7.364296	10.349618	21
CACCA	34480	7.121608	9.899373	40
AAACC	33430	6.9047375	9.918164	2
AGACG	3341590	6.7419705	86.54351	27
ACCCA	31570	6.520567	10.190249	33
ACCAC	31020	6.406968	9.414038	45
CAACC	29970	6.1900983	9.316799	31
CGGGT	12198075	5.9941854	235.08815	1
CCCAA	28735	5.935017	8.492127	41
CCAAC	28560	5.898872	9.607948	30
TGGCG	1182555	5.8796134	29.313614	30
CGGAG	5053650	5.822808	176.0491	1
CGGTT	15025490	5.514238	200.7726	1
CGGGG	7966295	5.241762	126.74163	1
GAGAC	2585430	5.216347	83.00835	26
AGGCG	4521660	5.2098494	96.63218	47
ACGCG	429870	5.0113297	19.30352	4

CGCGT	1000475	4.974319	27.562624	31
TCGAG	5607215	4.824946	85.34573	44
CGTCG	907590	4.512499	16.917011	41
CTCCA	50960	4.4890223	105.52665	24
GGGCG	6788470	4.466762	118.954834	2
CGAGG	3825075	4.407245	108.15149	45
CGGAC	375835	4.381402	145.79779	1
CGGTC	844050	4.196581	155.97433	1
TTACG	6346145	4.0782437	73.088264	14
GACGG	3379060	3.8933477	48.629208	28
CGCGA	332655	3.8780186	21.587734	5
AGCAC	189785	3.8742187	27.039764	10
CACCT	43235	3.808534	100.97362	31
CGGTA	4410195	3.794924	133.70163	1
TCACC	42670	3.758763	102.21534	30
ACGGG	3246030	3.7400706	48.660877	29
TACGT	5781530	3.7154033	73.57366	15
ACTCC	41875	3.6887321	105.25782	23
ACGTT	5723945	3.6783974	74.98909	16
AAAAA	3431830	3.67443	8.147603	31
GTCGA	4039310	3.4757807	84.87427	43
CGACG	292090	3.4051208	17.49039	24
GCACA	164225	3.3524437	26.982191	11
CGTTT	12225820	3.350833	44.580654	17
GAGGC	2779570	3.2026165	78.69273	46
ACACA	88580	3.1663892	3.888448	40
AGAGA	9062165	3.164342	27.717123	25
CGGAT	3653850	3.1440976	109.108	1
CCAGC	26465	3.121593	4.7110453	47
GGCGG	4725010	3.109021	25.292166	11
GGCGT	5932750	2.9153783	49.089783	3
CGGTG	5821955	2.8609333	51.84173	7
CGAGA	1364225	2.7524517	33.974266	25
ATCGC	312600	2.7215903	42.621006	29
GGAGG	22981745	2.6171088	32.5346	39
TTTCG	9505005	2.6051161	11.870341	30
AGGTC	2966550	2.552683	81.01063	41
ACGGA	1261930	2.5460618	15.791202	30
GGTCG	5136830	2.52426	49.736893	42
GCGGC	372260	2.478311	7.806432	6
AAGGC	1228005	2.4776149	41.516808	8
TTTTT	163052245	2.463479	5.2058516	16
TTCGA	3779450	2.4288003	24.943624	31
GGGAG	21085555	2.401175	30.191902	38
AGCGA	1173115	2.3668692	42.16952	9
GCGGT	4690755	2.3050568	42.721626	6
TTTAC	4786130	2.2970245	53.52729	13
TTTTA	64677370	2.2911994	12.757699	26
CGTTA	3499775	2.2490718	19.557167	9
AGATC	1485190	2.2378635	17.364532	27
CGTAG	2598670	2.2361264	27.537935	5
TCGTT	8150355	2.2338362	5.3644037	37
TTTAG	46988775	2.2288804	16.290613	26
AGTAG	14825170	2.2078154	14.685189	35
GCGGG	3338915	2.196981	25.663765	12
ATTTT	61411330	2.1755	10.123474	25
ACACG	106425	2.1725304	25.864697	13
GAGGT	25498345	2.1685464	24.67359	40
GACGC	184365	2.1492867	19.388512	3
CGTTC	574850	2.1345184	22.502752	33
GCGGA	1848665	2.1300287	21.452438	7
TTTCG	7760615	2.1270168	5.337359	35
TTCGC	562830	2.0898862	5.127784	13
ATCGT	3200320	2.056632	20.473682	39
ACGTC	235930	2.054078	11.730646	15
TAGTA	18452225	2.0522478	23.293158	29
AATTT	24687580	2.0505834	21.511566	24
TAGAG	13718655	2.0430293	12.7639	24
TACGG	2308500	1.9864384	20.497635	5
GGAA	13268090	1.9759297	12.64153	2
ACGGC	168815	1.968008	13.067123	12
ATTCC	3009810	1.934204	26.981544	34
AGGAG	9691705	1.9326149	6.7070365	38
AACGG	951205	1.9191451	14.802889	39
TATCG	2980040	1.9150728	21.000946	28
AGCGC	163920	1.9109434	9.91743	35
GGAAG	9507300	1.8958429	11.668516	2
CGAGT	2189425	1.8839757	29.91091	33
AAACG	529440	1.8704959	5.5637145	7
GCGTT	5034755	1.8477159	16.93418	16
TGGAA	12304100	1.832369	8.867702	1
TTAGT	38615085	1.83168	15.651641	28
GTCGC	367425	1.8268214	9.578081	3
GTAGA	12262295	1.8261433	12.470272	23
487005		1.8083346	7.7660403	40
GCCTA	2101375	1.8082099	27.2074	4
GAACG	893200	1.8021146	14.379509	28
TACCC	202805	1.7656817	7.5661135	13
GATCC	2010240	1.729789	11.003488	28
TAATT	20592530	1.7104434	21.304878	23
GAAA	4862980	1.6980639	11.736939	2
TAITTT	47354510	1.6775361	8.5669	32
GGAGA	8334835	1.6620426	10.598725	2
GGACG	1440120	1.6593039	15.097048	2
TAGTT	34719945	1.6469167	6.0624456	29
GACCC	1426595	1.6437205	8.4733305	37
TGGGA	19318110	1.6429385	19.5347	28
GAGAT	11006820	1.6391736	10.212023	26
AACCC	79970	1.6324854	10.531976	11
GTACC	1895935	1.6314309	20.3167	4
ACGTT	1889595	1.6259754	19.942123	6
AGGTA	10792515	1.6072586	20.78472	47
AGTCC	1816005	1.562652	14.423888	22
CGTGG	3171275	1.558378	28.939718	5
ACGAG	764745	1.5429446	5.651633	32
GGTTT	56746785	1.5371897	8.9232	2
AGTTT	32200635	1.527415	6.130693	26

TCGAA	1006490	1.5165651	5.1363153	44
AGTTA	13490640	1.5004226	11.836399	30
AGCGT	1740400	1.497595	6.9270887	29
CGATT	2313435	1.4866903	15.437219	11
GTTAA	13347040	1.4844514	27.15877	3
GGTTA	23343000	1.482625	21.391153	2
GGGAA	7385675	1.4727714	13.005954	2
AGGTT	23187395	1.4727417	13.109408	41
GTCGT	4005100	1.4698404	9.04289	3
GCGTG	2979110	1.4639473	29.031385	4
ATGGA	9786895	1.4574981	6.0017304	10
TAGGA	9783940	1.4570581	5.1608195	37
C CGTC	292840	1.4559882	9.414934	4
CGTAC	164450	1.4317516	7.6929317	13
GAATG	9515185	1.4170341	5.956864	13
GTTTA	29679000	1.4078028	11.208094	12
AATGG	9404990	1.4006233	6.008116	19
GGCGA	1213965	1.3987285	9.354019	2
AGATA	5265630	1.3731556	6.355313	26
TTGTA	28937410	1.3726261	15.167635	20
GTAAT	12329855	1.3713205	26.96584	22
TCGTG	3719525	1.3650366	11.132434	40
TTTTG	66908180	1.3535783	5.5735526	34
CGAAG	668340	1.3484385	5.621409	45
TTAAG	12096770	1.345397	7.594661	6
GCGAC	115375	1.3450165	11.412874	23
TTTAA	16122095	1.339123	6.2845526	5
TGGCG	2681570	1.3177347	22.00641	10
GGAGT	15479490	1.3164773	10.435777	2
GTAGT	20574890	1.3068092	6.735604	36
GGTAG	15360710	1.3063753	7.565421	2
AAGGC	642575	1.296455	26.979368	46
GGGTT	35543445	1.2892222	12.608119	2
TTCCG	3507370	1.2871773	15.214505	35
TTATA	15493410	1.2869036	8.733938	46
AGTAT	11540875	1.2835705	22.384232	30
TATAG	11532355	1.282623	11.2535	47
AAGTA	4888040	1.2746888	7.362414	34
GCGAT	1475910	1.2700043	18.948704	10
GGAA	623585	1.2581412	13.518322	27
GACGT	1461055	1.2572217	5.7469044	3
TGTAA	11236145	1.2496786	26.651026	21
GTATT	26149285	1.2403733	10.044231	31
AAAAC	198695	1.2292305	7.886452	6
ATTAT	14514440	1.2055889	8.561851	45
TTTGT	59275200	1.1991601	7.4399195	19
TTAAT	14339990	1.1910989	18.748163	4
GGGGA	10404925	1.184889	10.544952	2
CGTGT	3226270	1.184016	10.67358	41
TTGAG	18636055	1.1836646	10.380165	44
GGTGG	24248080	1.1776811	12.09974	8
GAGTA	7720230	1.1497233	10.560885	34
GGGGT	23443115	1.1385854	8.604351	2
TCGGA	1308030	1.1255454	5.3003993	46
GGGAT	13154635	1.1187563	9.393356	2
TCGAT	1724460	1.1081953	5.196117	11
TTATC	2298105	1.1029378	15.406057	37
TAAGC	726080	1.0940473	29.843534	7
CGTAA	722310	1.0883666	6.048406	21
GGGTA	12741050	1.0835824	15.4232235	2
TAGGC	1257745	1.0822755	5.759369	41
TCGGG	2178710	1.0706273	19.233263	36
GATTA	9625360	1.070528	10.930311	44
TGGAG	12528780	1.0655295	9.374365	1
TGAGG	12406505	1.0551304	12.762529	45
CGAAC	51540	1.0521233	8.479301	9
GGATT	16500675	1.0480363	6.439781	2
AAGGT	6995120	1.0417373	5.1407375	46
AGTAA	3923085	1.0230507	5.1945643	9
ATTAC	896995	1.0093913	6.6427464	29
TGCGG	2026320	0.99574214	9.9241905	5
TTTTT	4861395	0.99507	8.123514	29
TTGGG	26967020	0.9781404	8.857491	36
TAAGG	6558555	0.9767226	7.2268596	45
GGTAC	1122215	0.9656535	20.469666	3
TGGTT	35338200	0.95726156	7.013125	1
TTTTG	35226545	0.95423704	7.0833845	35
GAATA	3638910	0.94894433	5.1759977	3
GTTAT	19891650	0.94354665	5.2401175	31
TGGGG	19354725	0.9400204	8.664409	1
GGTTG	25866265	0.9382142	5.901542	42
TGTAG	14741005	0.93627137	5.18033	21
CGTGA	1084495	0.93319595	6.5819983	26
GGATA	6197225	0.9229121	7.612534	2
GTGGT	25371430	0.9202656	10.035575	9
AGCTG	10806250	0.9190342	5.621419	47
GTTTG	33454965	0.9062474	7.574434	18
AGTTG	14179955	0.9006364	6.7021933	38
TATTC	1869195	0.8970895	18.011375	33
TAGAC	580745	0.8750585	6.886058	25
GGAGC	759425	0.87500834	7.5746193	27
GTTGA	13564515	0.8615468	10.114005	43
GGCTG	17637145	0.8566011	8.888792	2
AGTGA	5550130	0.82654446	5.584361	18
GGTAT	12778530	0.81162524	6.3113537	2
GTCCG	1651215	0.81141394	9.2316095	4
CGATC	92665	0.80676955	6.298049	44
GGGGG	12187695	0.79259974	6.1551924	2
AGTGG	9292480	0.7902933	6.5669203	8
GGTAA	5197050	0.7739626	6.481721	2
GTCCG	1572195	0.7725832	20.672384	9
CACGT	88505	0.7705513	11.047466	14
TGGGT	20787710	0.75400615	8.633229	1
GAAGC	372790	0.7521386	5.0776057	4
TGCTG	19375595	0.70278627	5.4425793	1
AACTC	45985	0.7010619	18.976315	22
GAGTC	803065	0.6910285	13.164964	21

TGGTA	9908115	0.62931144	5.1151824	1
AAGTC	378650	0.57054454	5.301346	41
GGTGC	942705	0.46324924	9.33956	3
TCCAG	51690	0.45002884	10.814286	25
CAGTC	46620	0.40588784	10.988149	27
GTCAC	44140	0.3842962	10.849058	29
CCAGT	43015	0.37450162	10.707923	26
ATGCC	42450	0.3695826	11.55703	47
ATCTC	50195	0.3263715	8.663195	40
ACCTT	41355	0.26889318	7.662386	32



## 5 Overrepresented sequences

Note: FastQC tracks sequences that appear in the first 200,000 reads to the end of the file.

Sequence	Count	%	Possible Source
CGGGCGCGGTGGTTACGTTTGTAATTTAGTATTTTGGGAGGTCGAGGCG	277819	0.34421462001292463	No Hit
CGGTTAATTTTGTATTTTAGTAGAGACGGGGTTTATCGTGTAGTTA	150587	0.18657560132275433	No Hit
CGGGTTACGTTATTTTGTGTTTAGTTTTTCGAGTAGTTGGGATTATAG	149553	0.1852944869385928	No Hit
CGGTTAATTTTGTATTTTAGTAGAGACGGGGTTTATTTTGTAGTTA	97054	0.12024881570639297	No Hit
CGGGTTACGTTATTTTGTGTTTAGTTTTTAAGTAGTTGGGATTATAG	89386	0.1107482498478336	No Hit