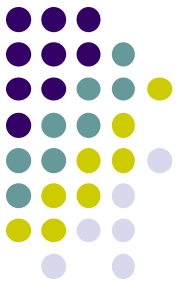




FastQC

Quality Assurance tool for FASTQ sequences

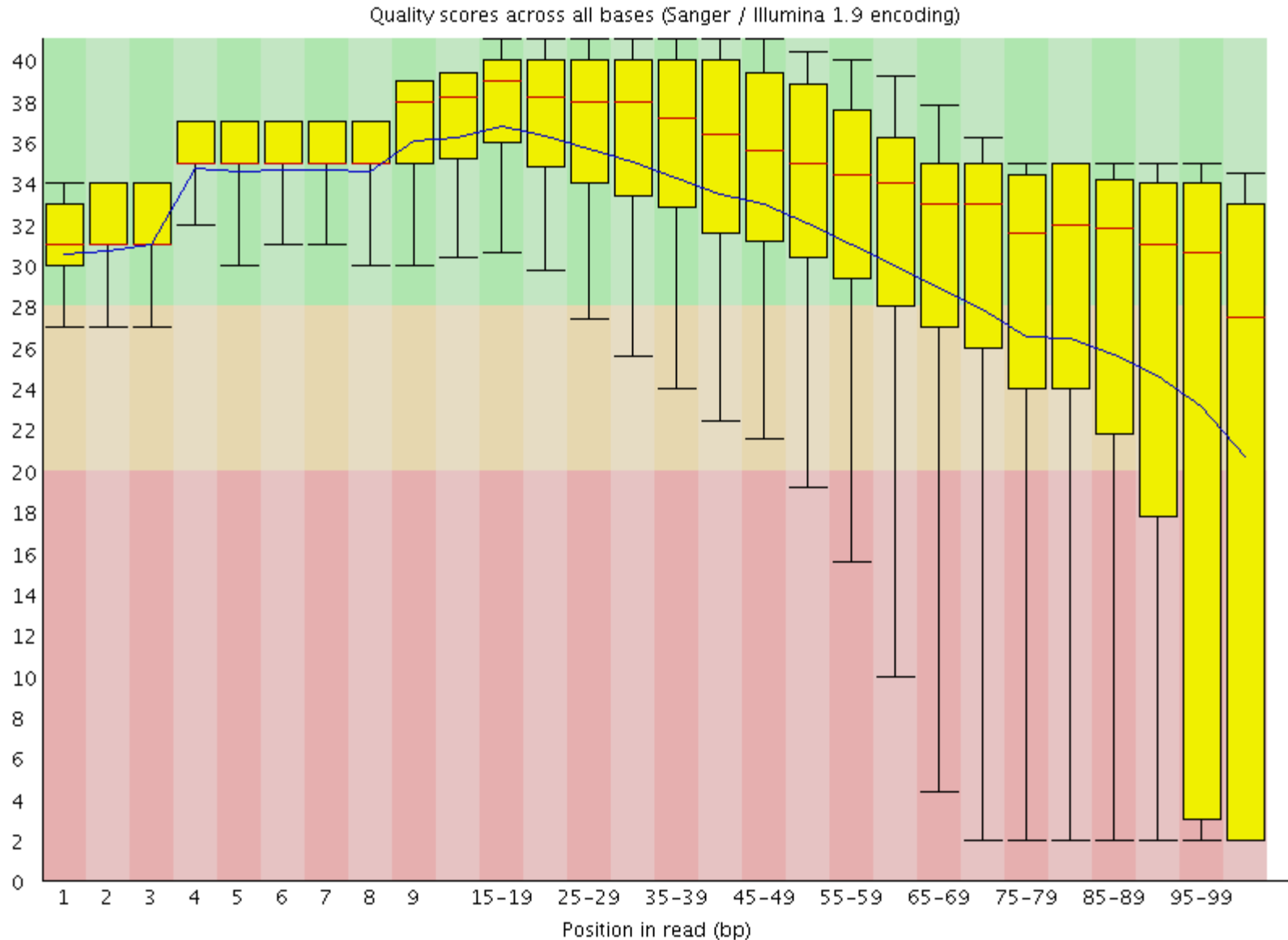
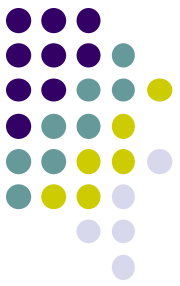
- FastQC website:
<http://www.bioinformatics.babraham.ac.uk>
- FastQC report documentation:
<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/Help/3%20Analysis%20Modules/>
- Good Illumina dataset:
http://www.bioinformatics.babraham.ac.uk/projects/fastqc/good_sequence_short_fastqc/fastqc_report.html
- Bad Illumina dataset:
http://www.bioinformatics.babraham.ac.uk/projects/fastqc/bad_sequence_fastqc/fastqc_report.html
- Real Yeast ChIP-seq dataset:
http://web.corral.tacc.utexas.edu/BiolTeam/yeast_stuff/Sample_Yeast_L005_R1.cat_fastqc/fastqc_report.html



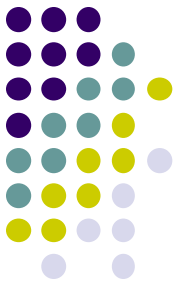
Most useful reports

- Should I trim low quality bases?
 - Per-base sequence quality Report
 - based on *all* sequences
- Do I need to remove adapter sequences?
 - Overrepresented sequences Report
 - based on *1st 200,000* sequences
- How complex is my library?
 - Sequence duplication levels Report
 - estimate based on *1st 200,000* sequences

Per-base sequence quality



Overrepresented Sequences



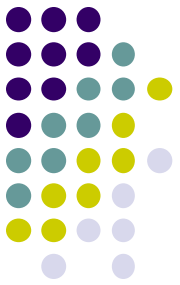
| Sequence | Count | Percentage | Possible Source |
|--|-------|---------------------|--|
| AGATCGGAAGAGCACACGTCTGAACTCCAGTCACCTCAGAATCTCGTATG | 60030 | 5.01369306977828 | TruSeq Adapter, Index 1 (97% over 37bp) |
| GATCGGAAGAGCACACGTCTGAACTCCAGTCACCTCAGAATCTCGTATGC | 42955 | 3.5875926338884896 | TruSeq Adapter, Index 1 (97% over 37bp) |
| CACACGTCTGAACTCCAGTCACCTCAGAATCTCGTATGCCGTCTTCTGCT | 3574 | 0.29849973398946483 | RNA PCR Primer, Index 40 (100% over 41bp) |
| CAGATCGGAAGAGCACACGTCTGAACTCCAGTCACCTCAGAATCTCGTAT | 2519 | 0.2103863542024236 | TruSeq Adapter, Index 1 (97% over 37bp) |
| GAGATCGGAAGAGCACACGTCTGAACTCCAGTCACCTCAGAATCTCGTAT | 1251 | 0.10448325887543942 | TruSeq Adapter, Index 1 (97% over 37bp) |

Overrepresented Sequences



| Sequence | Count | Percentage | Possible Source |
|--|--------|---------------------|---|
| AACGACTCTCGGCAACGGATATCTCGGCTCTCGCATCGATGAAGAACGTA | 102020 | 1.0707851766890004 | No Hit |
| AATTCTAGAGCTAATACGTGCAACAAACCCCGACTTATGGAAGGGACGCA | 89437 | 0.9387160737848865 | No Hit |
| AAAGGATTGGCTCTGAGGGCTGGGCTCGGGGGTCCCAGTTCGGAACCCGT | 89427 | 0.9386111154260659 | No Hit |
| TACCTGGTTGATCCTGCCAGTAGTCATATGCTTGTCTCAAAGATTAAGCC | 87604 | 0.9194772066130483 | No Hit |
| ATTGGCTCTGAGGGCTGGGCTCGGGGGTCCCAGTTCGGAACCCGTGGCT | 65829 | 0.6909303802809273 | No Hit |
| TCTAGAGCTAATACGTGCAACAAACCCCGACTTATGGAAGGGACGCATT | 65212 | 0.6844544495416888 | No Hit |
| TAAAACGACTCTCGGCAACGGATATCTCGGCTCTCGCATCGATGAAGAAC | 61582 | 0.646354565289767 | No Hit |
| CTCGGATAACCGTAGTAATTCTAGAGCTAATACGTGCAACAAACCCCGAC | 59180 | 0.6211435675010296 | No Hit |
| ATGGATCCGTAACCTCGGGAAAAGGATTGGCTCTGAGGGCTGGGCTCGGG | 56982 | 0.598073720232235 | No Hit |
| AAAACGACTCTCGGCAACGGATATCTCGGCTCTCGCATCGATGAAGAACG | 54813 | 0.5753082522040206 | No Hit |
| ATTCTAGAGCTAATACGTGCAACAAACCCCGACTTATGGAAGGGACGCAT | 52688 | 0.5530046009546172 | No Hit |
| GCGACCCAGGTGAGGCGGGATTACCCGCTGAGTTTAAGCATATCAATAA | 41363 | 0.4341392595901502 | No Hit |
| CTAGAGCTAATACGTGCAACAAACCCCGACTTATGGAAGGGACGCATTTA | 40019 | 0.4200328561646452 | No Hit |
| AGAACTCCGCAGTTAAGCGTGCTTGGGCGAGAGTAGTACTAGGATGGGTG | 39753 | 0.4172409638200141 | No Hit |
| ACTCGGATAACCGTAGTAATTCTAGAGCTAATACGTGCAACAAACCCCGA | 38867 | 0.4079416532284981 | No Hit |
| ACGACTCTCGGCAACGGATATCTCGGCTCTCGCATCGATGAAGAACGTAG | 38438 | 0.40343893963508914 | No Hit |
| ACTTCGGGAAAAGGATTGGCTCTGAGGGCTGGGCTCGGGGGTCCCAGTTC | 37406 | 0.3926072370047907 | No Hit |
| AGATCGGAAGAGCACACGTCTGAACTCCAGTCACTGACCAATCTCGTATG | 34199 | 0.35894709133098535 | TruSeq Adapter, Index 4 (100% over 49bp) |
| GAACCTTGGGATGGGTGGCCGGTCCGCCTTTGGTGTGCATTGGTCGGCT | 34099 | 0.3578975077427782 | No Hit |

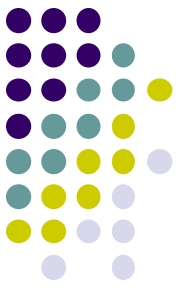
Overrepresented Sequences



| Sequence | Count | Percentage | Possible Source |
|---|---------|---------------------|-----------------|
| GAAGGTCACGGCGAGACGAGCCGTTTATCATTACGATAGGTGTCAAGTGG | 5632816 | 32.03026785752871 | No Hit |
| TATTCTGGTGTCTAGGCGTAGAGGAACAACACCAATCCATCCCGAACTT | 494014 | 2.8091456822607364 | No Hit |
| TCAAACGAGGAAAGGCTTACGGTGGATACCTAGGCACCCAGAGACGAGGA | 446641 | 2.539765344040083 | No Hit |
| TAAAACGACTCTCGGCAACGGATATCTCGGCTCTCGCATCGATGAAGAAC | 179252 | 1.0192929387357474 | No Hit |
| GAAGGTCACGGCGAGACGAGCCGTTTATCATTACGATAGGGGTCAAGTGG | 171681 | 0.9762414422996221 | No Hit |
| AACGACTCTCGGCAACGGATATCTCGGCTCTCGCATCGATGAAGAACGTA | 143415 | 0.8155105483274229 | No Hit |
| AGAACATGAAACCGTAAGCTCCCAAGCAGTGGGAGGAGCCCTGGGCTCTG | 111584 | 0.6345077504066322 | No Hit |
| AAAACGACTCTCGGCAACGGATATCTCGGCTCTCGCATCGATGAAGAACG | 111255 | 0.6326369351474214 | No Hit |
| ATTACGATAGGTGTCAAGTGGAAAGTGCAGTGATGTATGCAGCTGAGGCAT | 73682 | 0.41898300890326096 | No Hit |
| GAAGGTCACGGCGAGACGAGCCGTTTATCATTACGATAGGTGTCAAGGGG | 71661 | 0.4074908580252516 | No Hit |
| GGATGCGATCATACCAGCACTAATGCACCGGATCCCATCAGAACTCCGCA | 69548 | 0.3954755612388914 | No Hit |
| ATATTCTGGTGTCTAGGCGTAGAGGAACAACACCAATCCATCCCGAACT | 54017 | 0.30716057099328803 | No Hit |

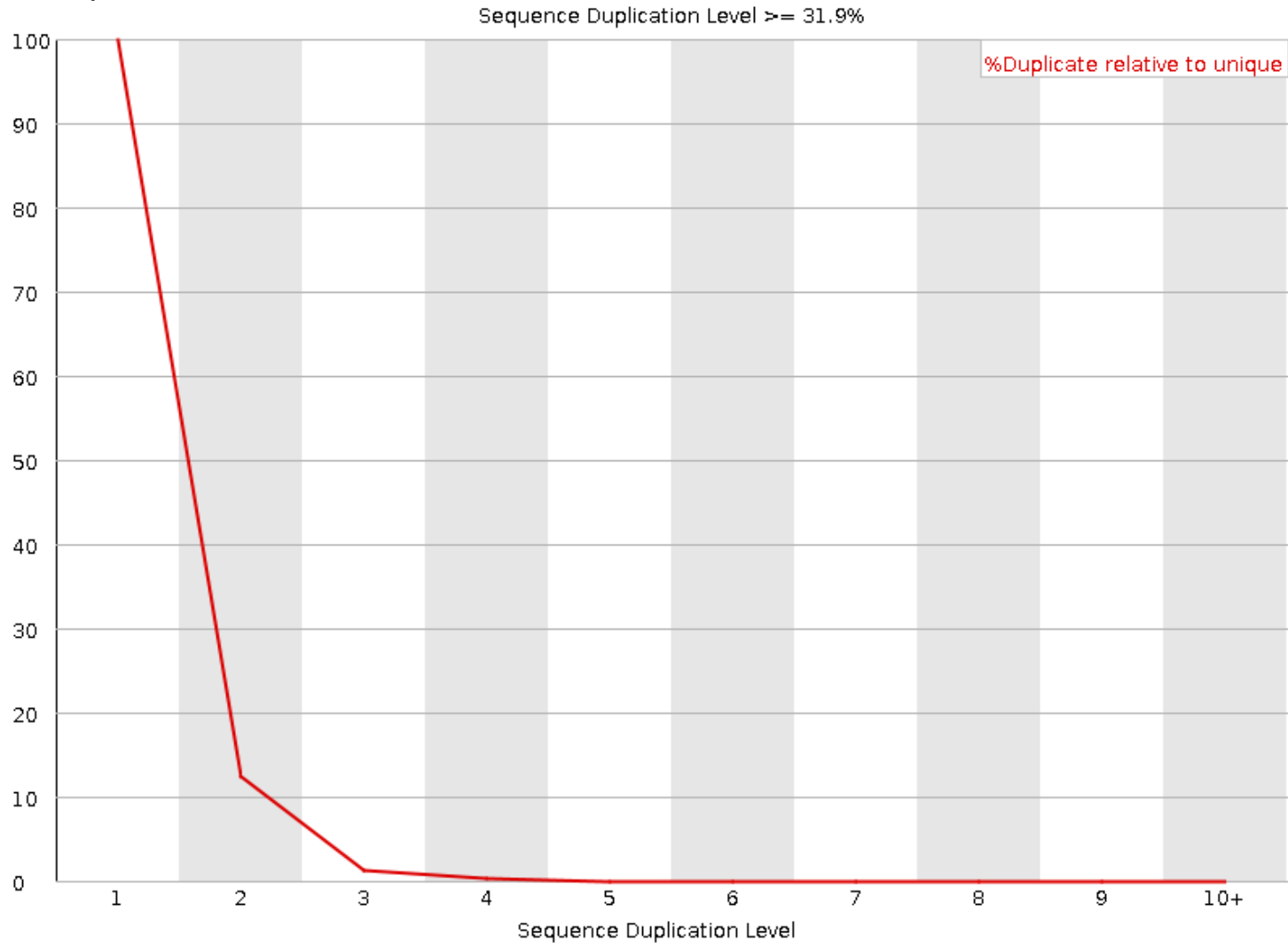
Duplication levels

Yeast ChIP-seq



for every 100 unique sequences
there are

- ~12 sequences w/2 copies
- ~1-2 with 3 copies

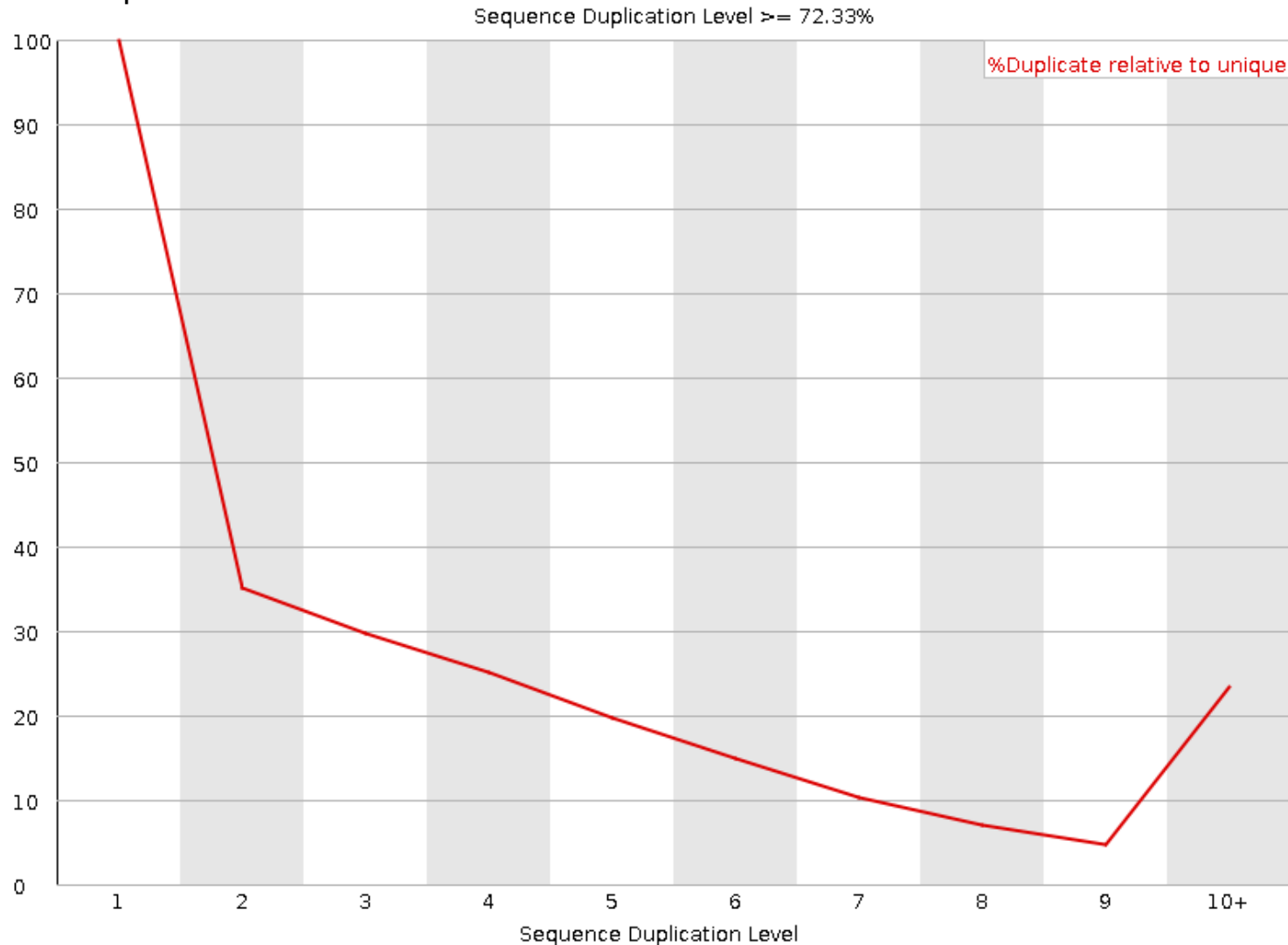
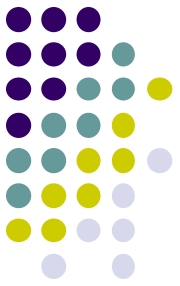


Duplication levels

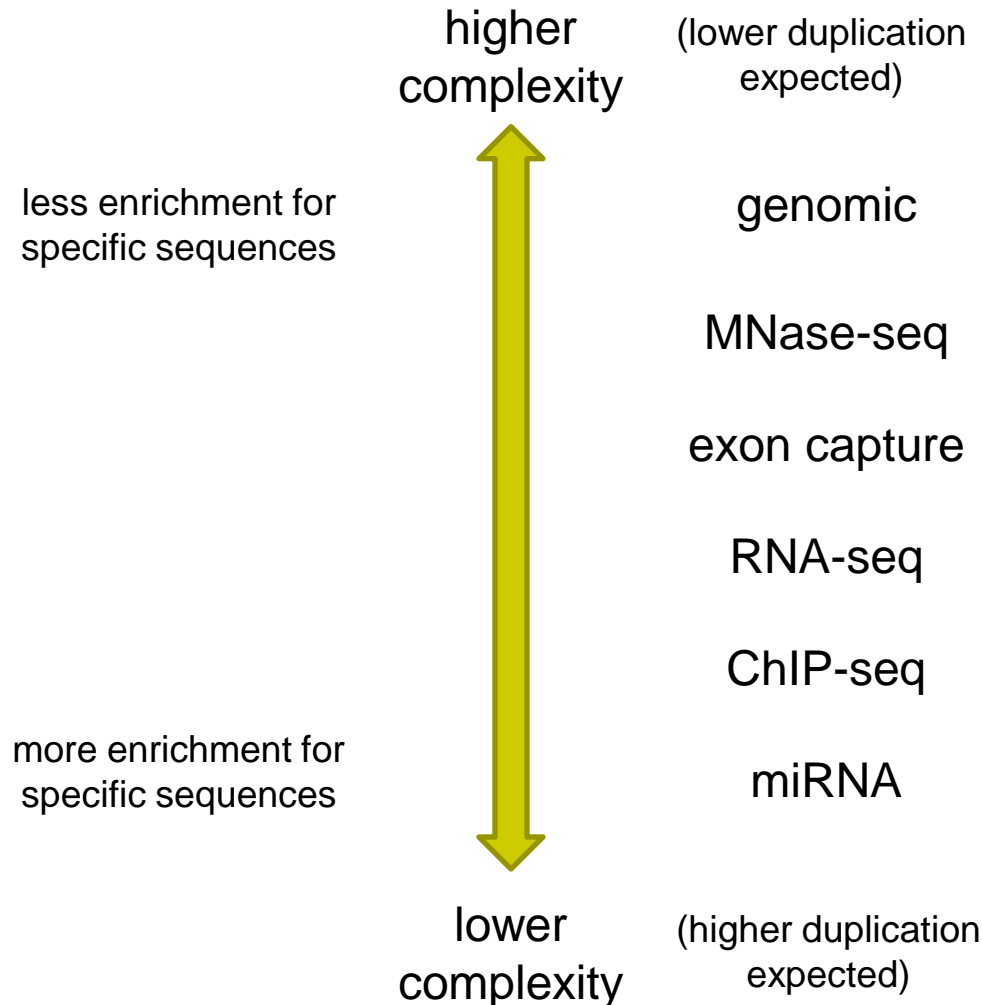
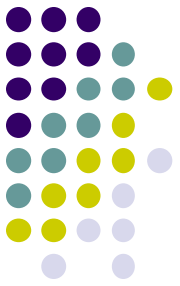
Yeast ChIP-exo

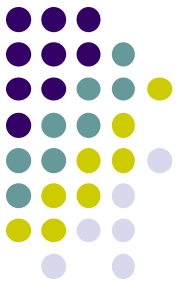
for every 100 unique sequences
there are

- ~35 sequences w/2 copies
- ~22 with 10+ copies



Library complexity is a function of experiment type (& sequencing depth)





Running FastQC

- Can run as interactive tool or command line
- Input:
 - fastq files (R1, R2 separately)
- Output:
 - directory with html & text reports
 - fastqc_report.html
 - fastqc_data.txt