

KBase Variation Services

Overview and Demo

Michael Schatz, James Gurtowski
Cold Spring Harbor Laboratory

1. Introduction to KBase
2. Resequencing and variation calling theory
3. KBase services for variation calling
4. Live Demo
5. Additional Resources

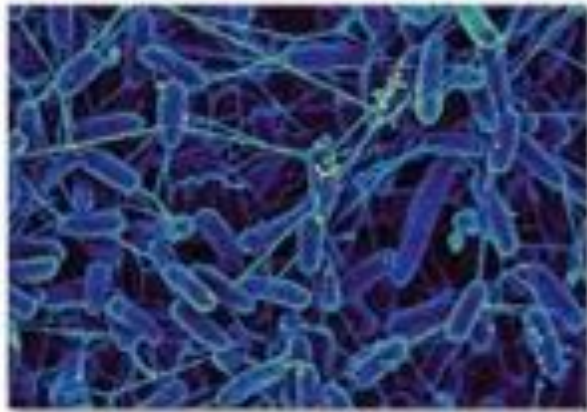


1. Introduction to KBase
2. Resequencing and variation calling theory
3. KBase services for variation calling
4. Live Demo
5. Additional Resources



Knowledgebase enabling ***predictive*** systems biology.

- Powerful ***modeling*** framework.
- ***Community-driven***, extensible and scalable ***open-source*** software and application system.
- Infrastructure for integration and reconciliation of ***algorithms*** and ***data sources***.
- Framework for standardization, search, and ***association*** of data
- Resources to enable ***experimental design*** and ***interpretation*** of results.



Microbes

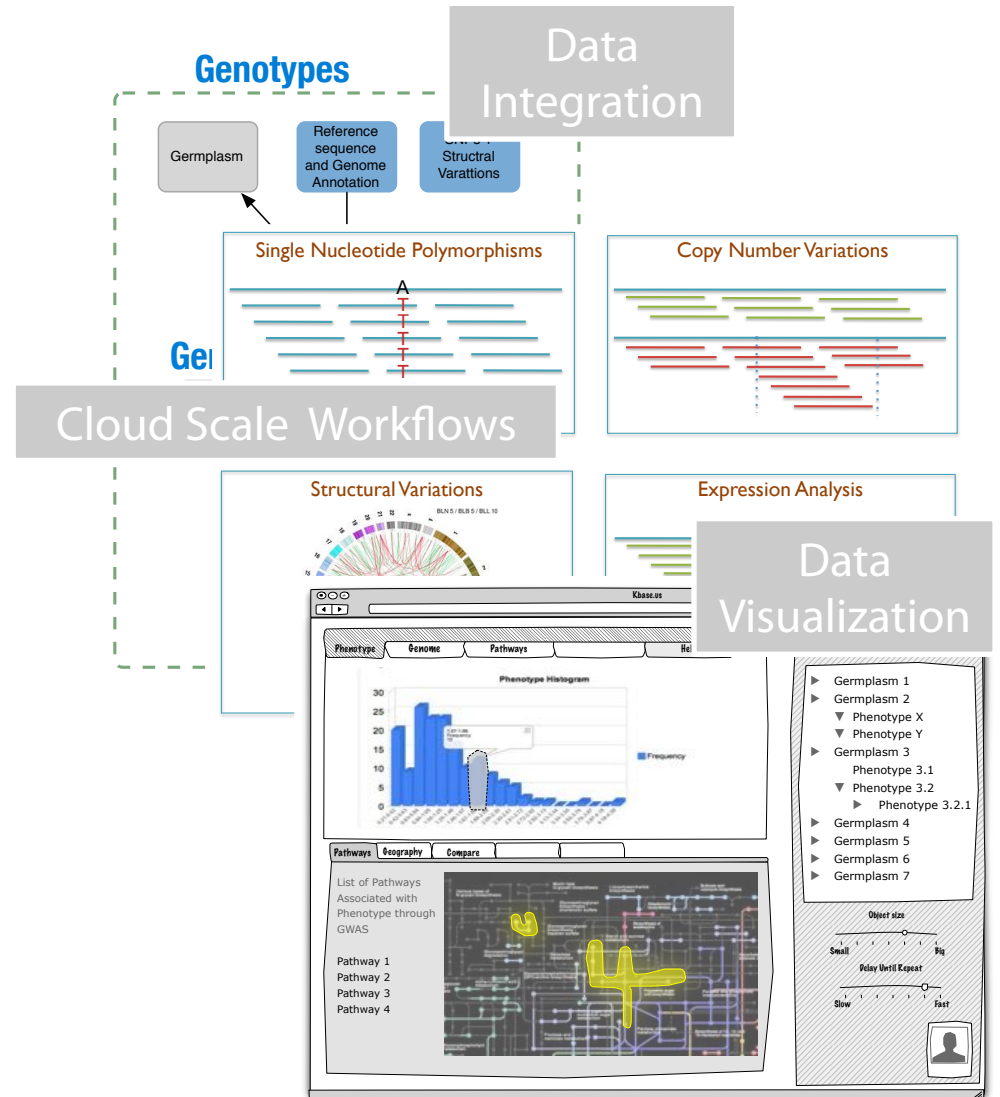
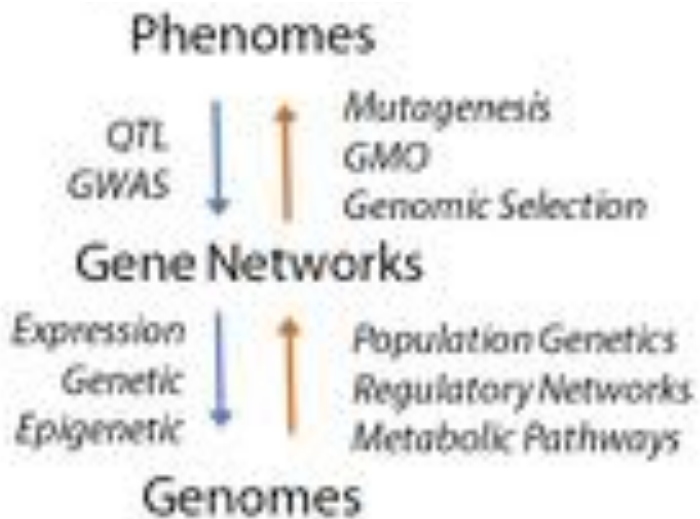


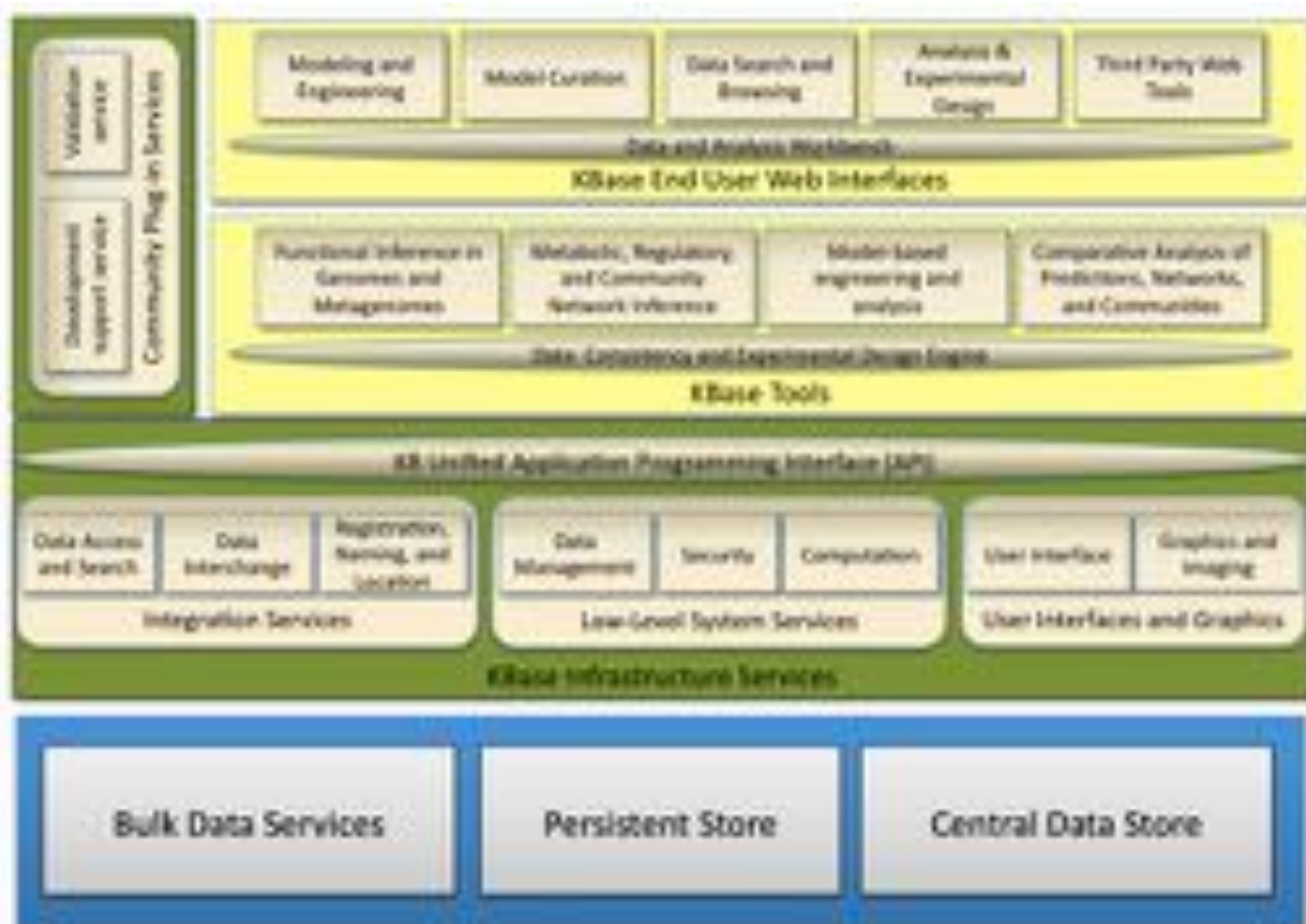
Communities



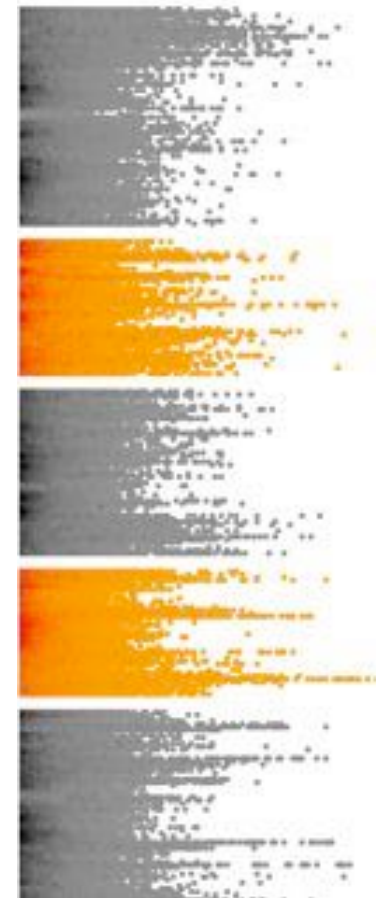
Plants

Model development
Hypothesis testing
Knowledge Synthesis





Variation Services: Samples to Discoveries

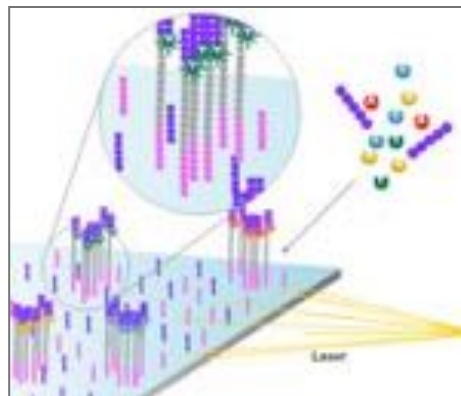
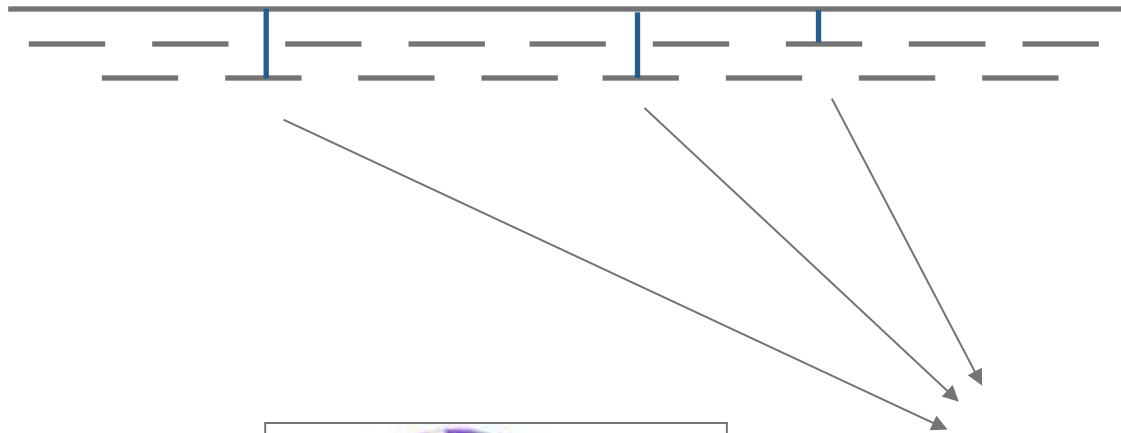


Powered by KBase

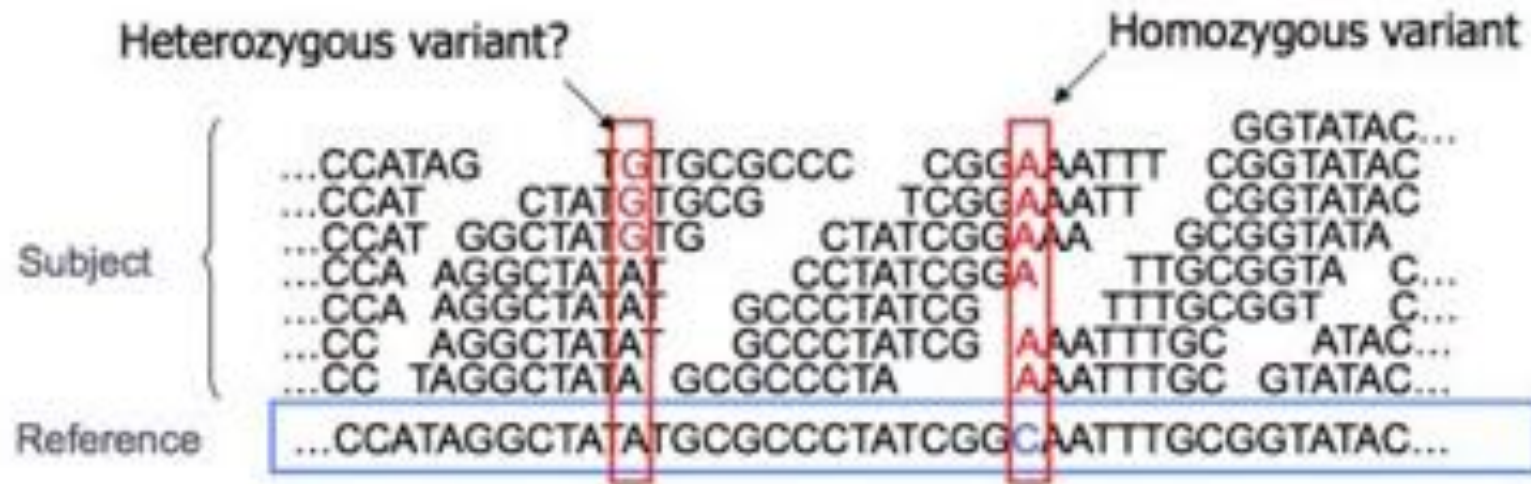
1. Introduction to KBase
2. Resequencing and variation calling theory
3. KBase services for variation calling
4. Live Demo
5. Additional Resources



How does your sample compare to the reference?



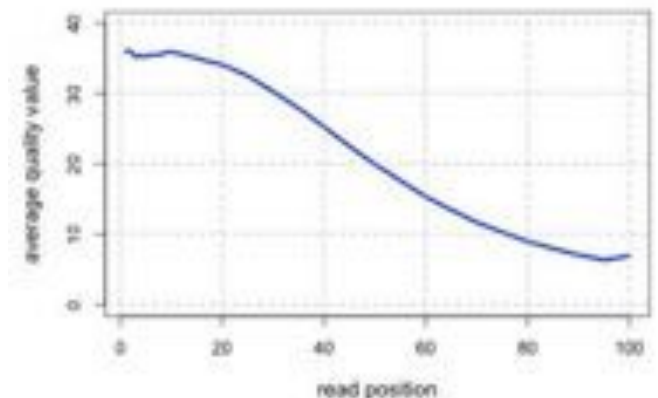
Plant Height —
 Drought Resistance —
 Biomass production —
 — — —
 — — —



- Sequencing instruments make mistakes
 - Quality of read decreases over the read length

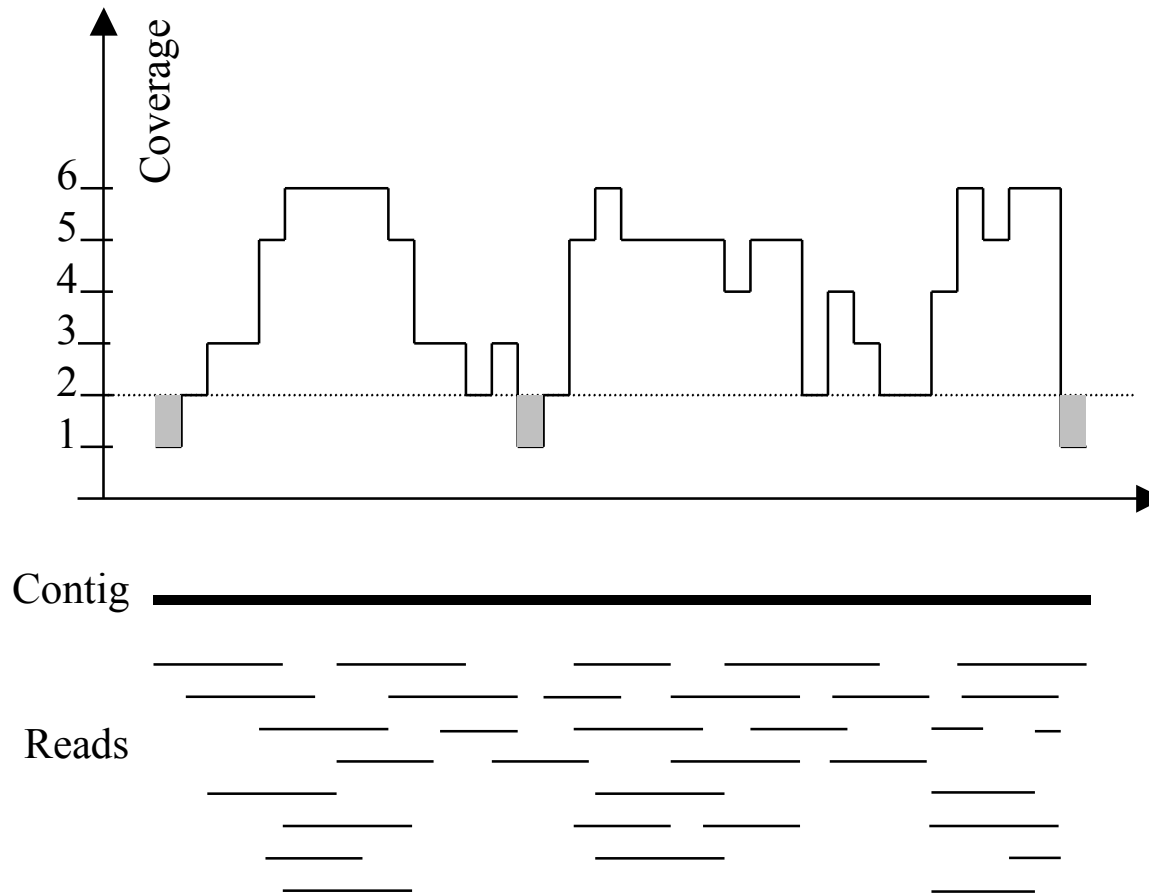
- A single read differing from the reference is probably just an error, but it becomes more likely to be real as we see it multiple times
 - Often framed as a Bayesian problem of more likely to be a real variant or chance occurrence of N errors
 - Accuracy improves with deeper coverage

$$Q_{\text{sanger}} = -10 \log_{10} p$$



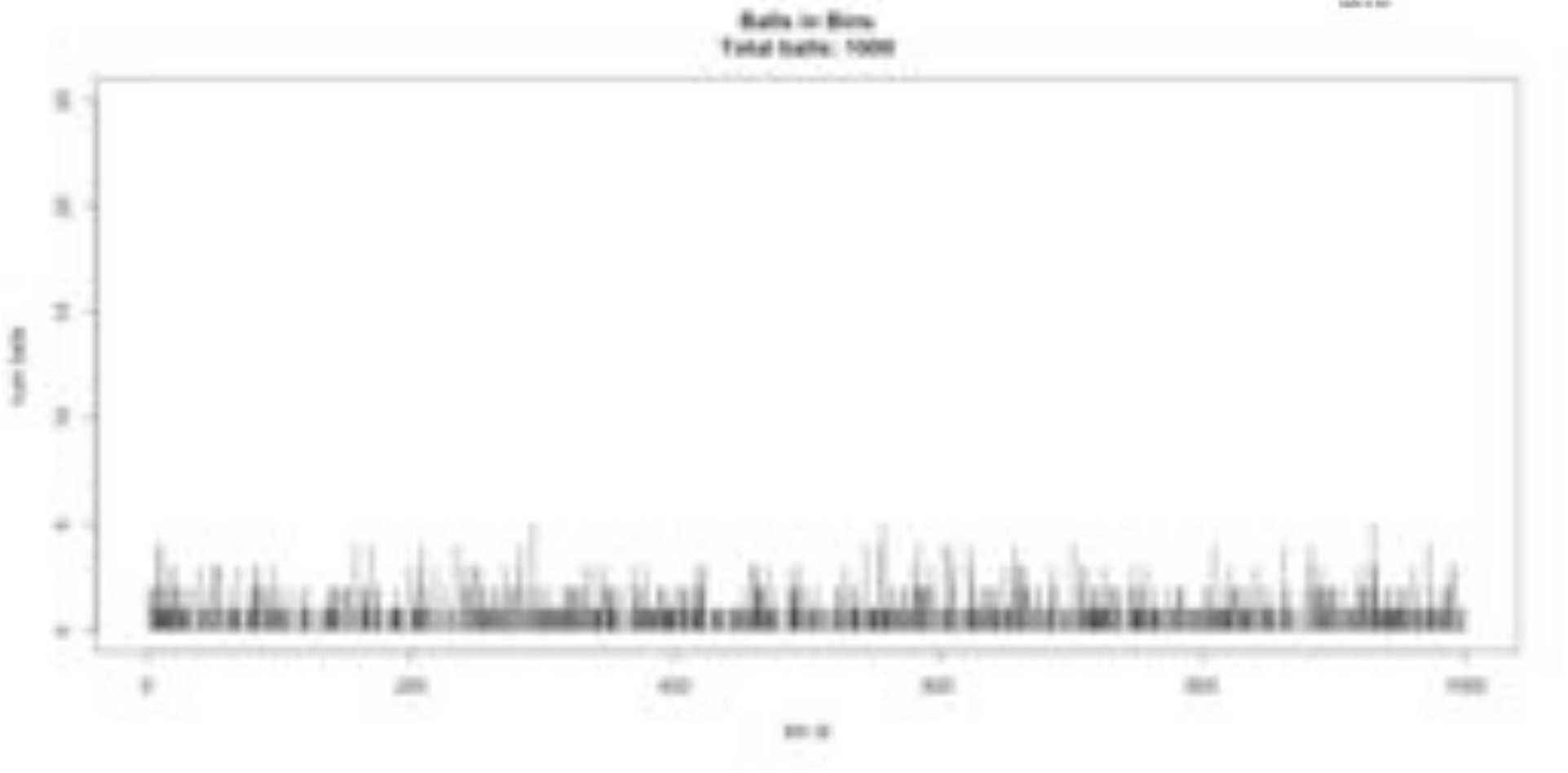
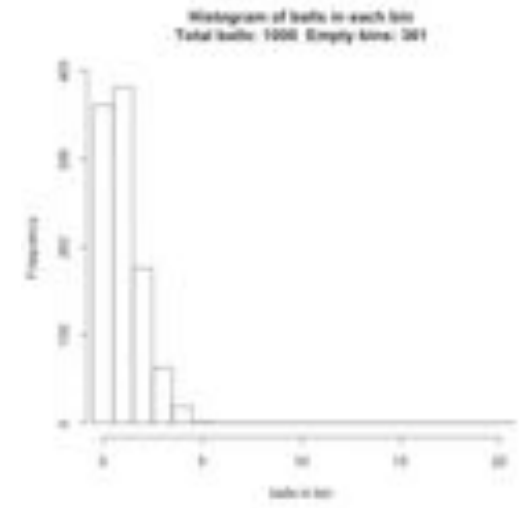
Coverage

Typical contig coverage

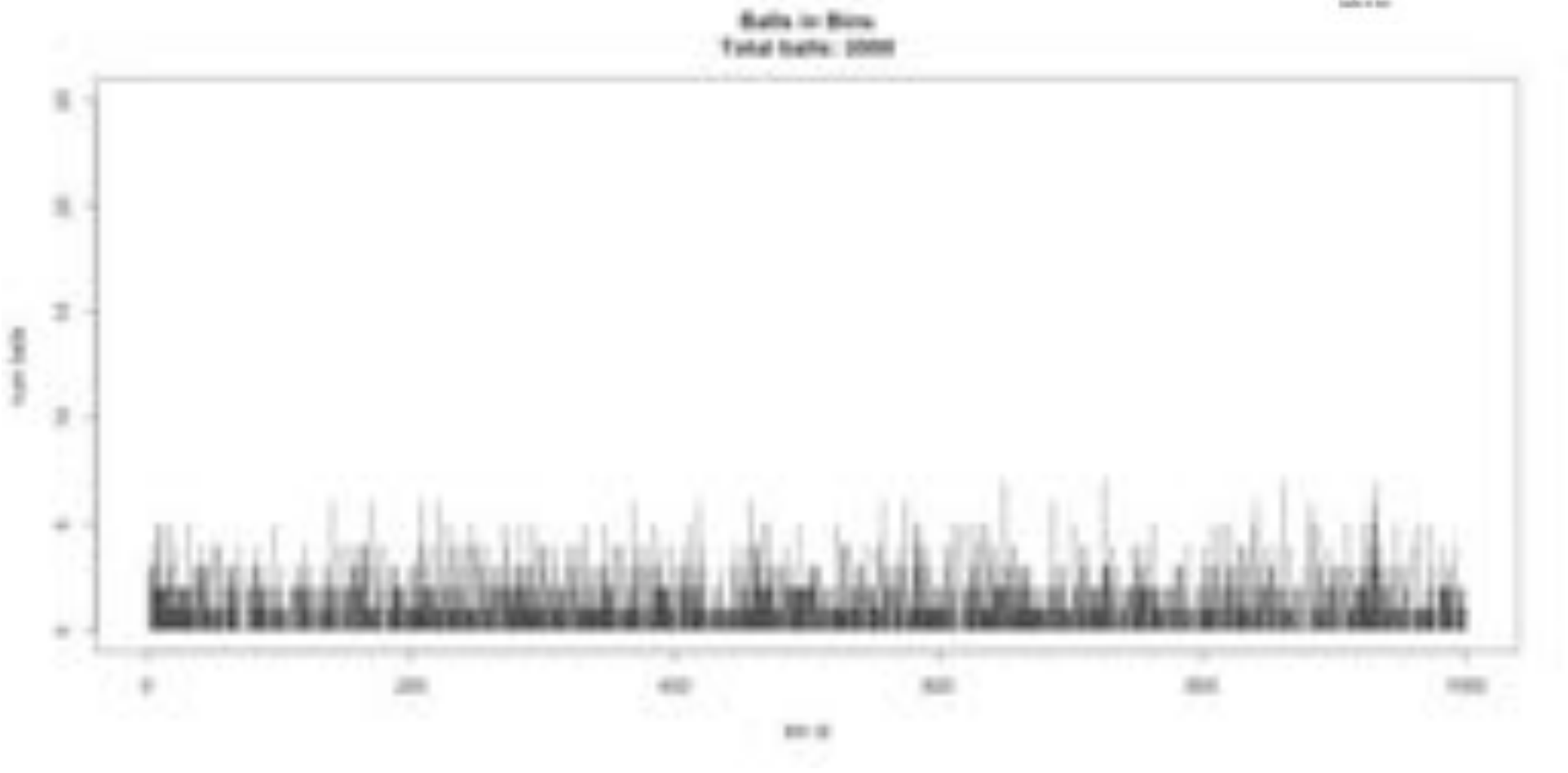
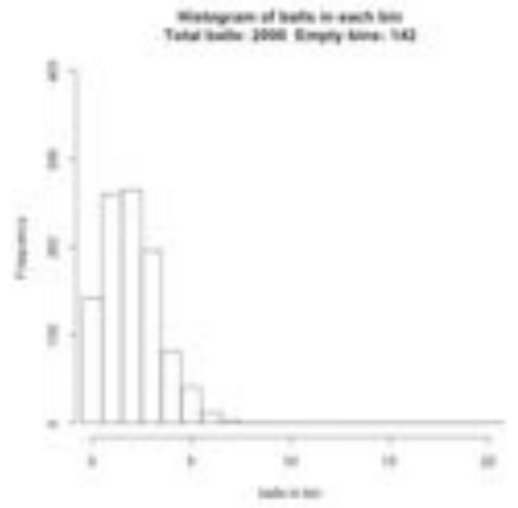


Imagine raindrops on a sidewalk

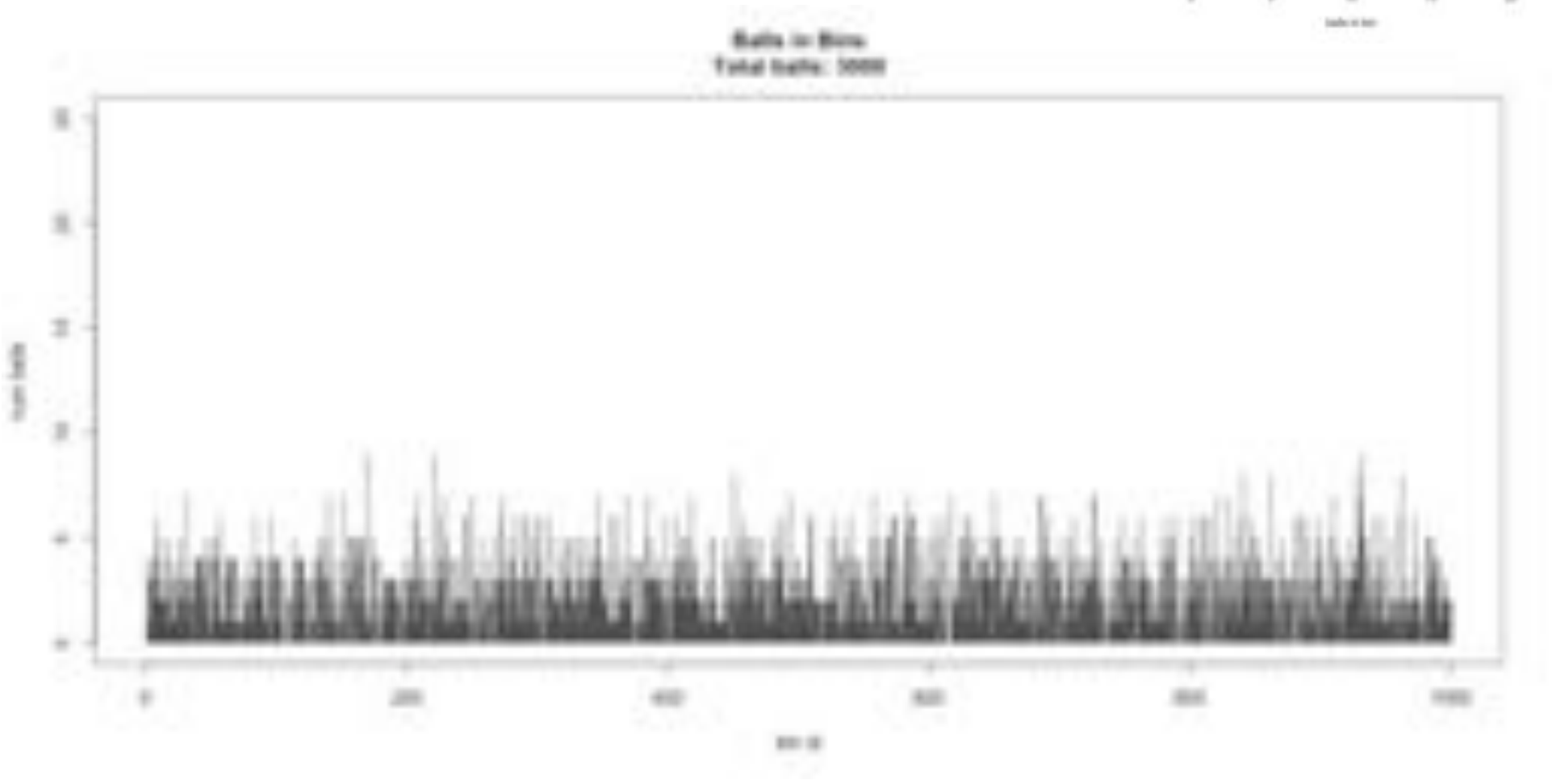
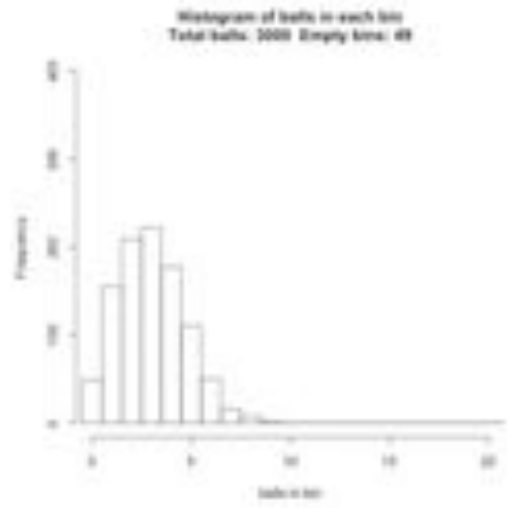
Ix Sequencing



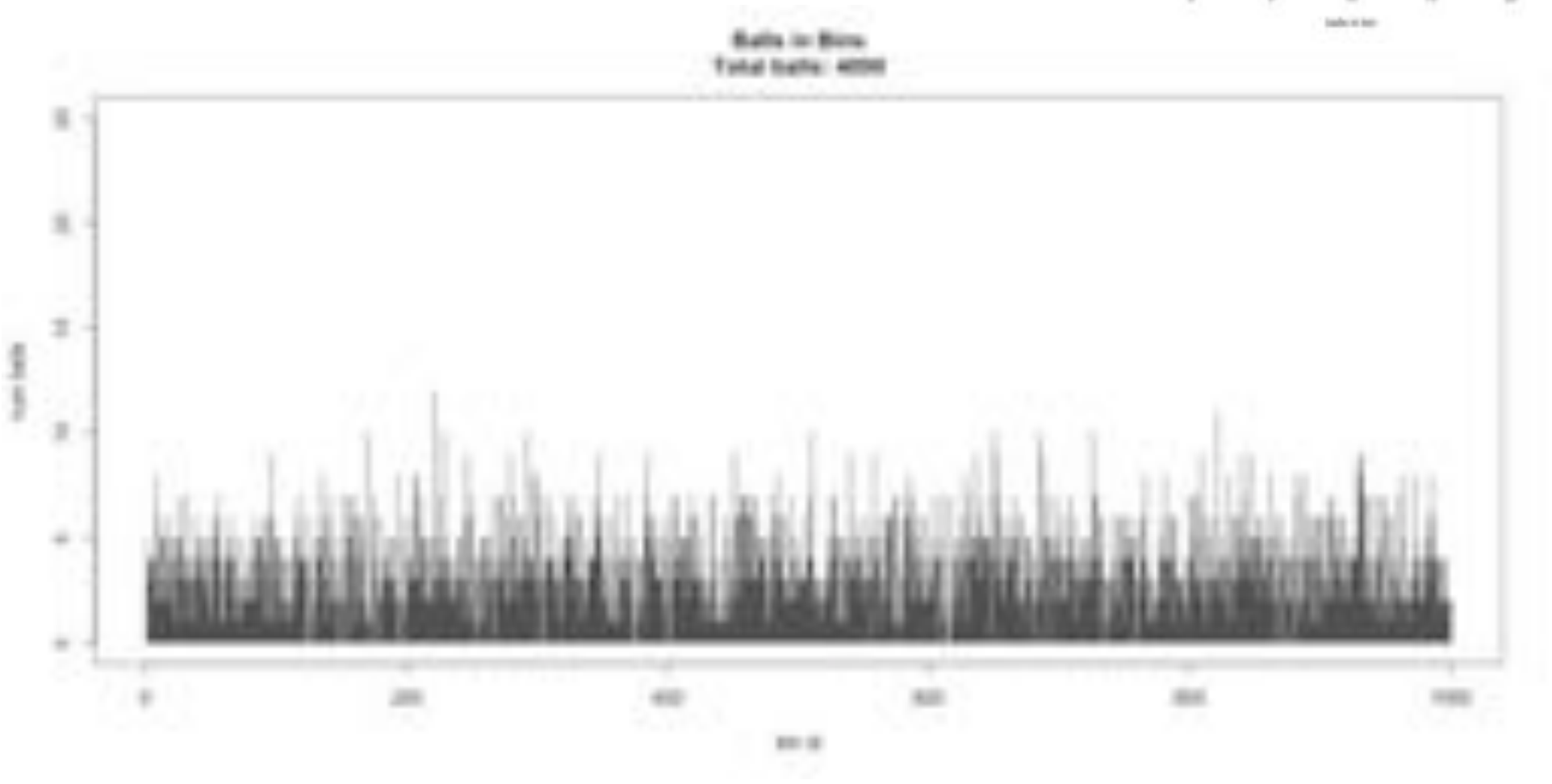
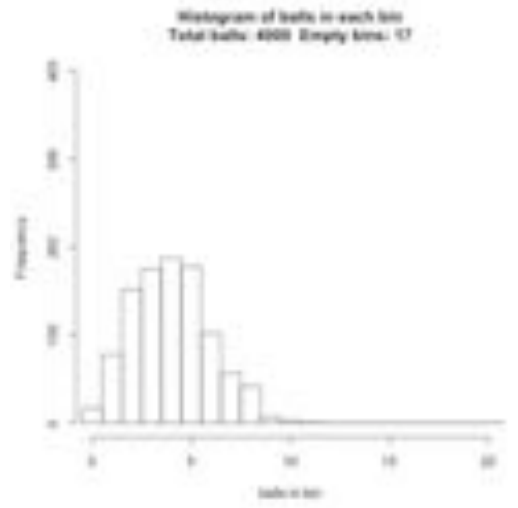
2x Sequencing



3x Sequencing

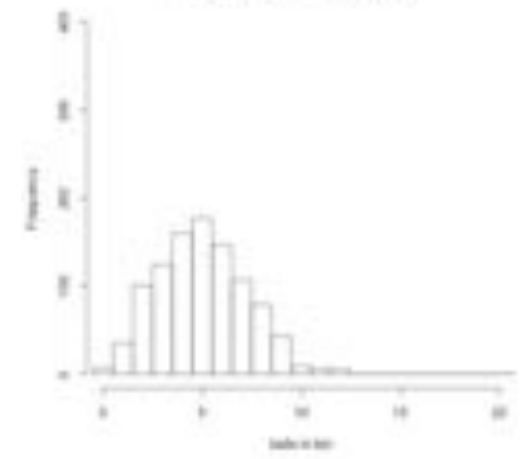


4x Sequencing

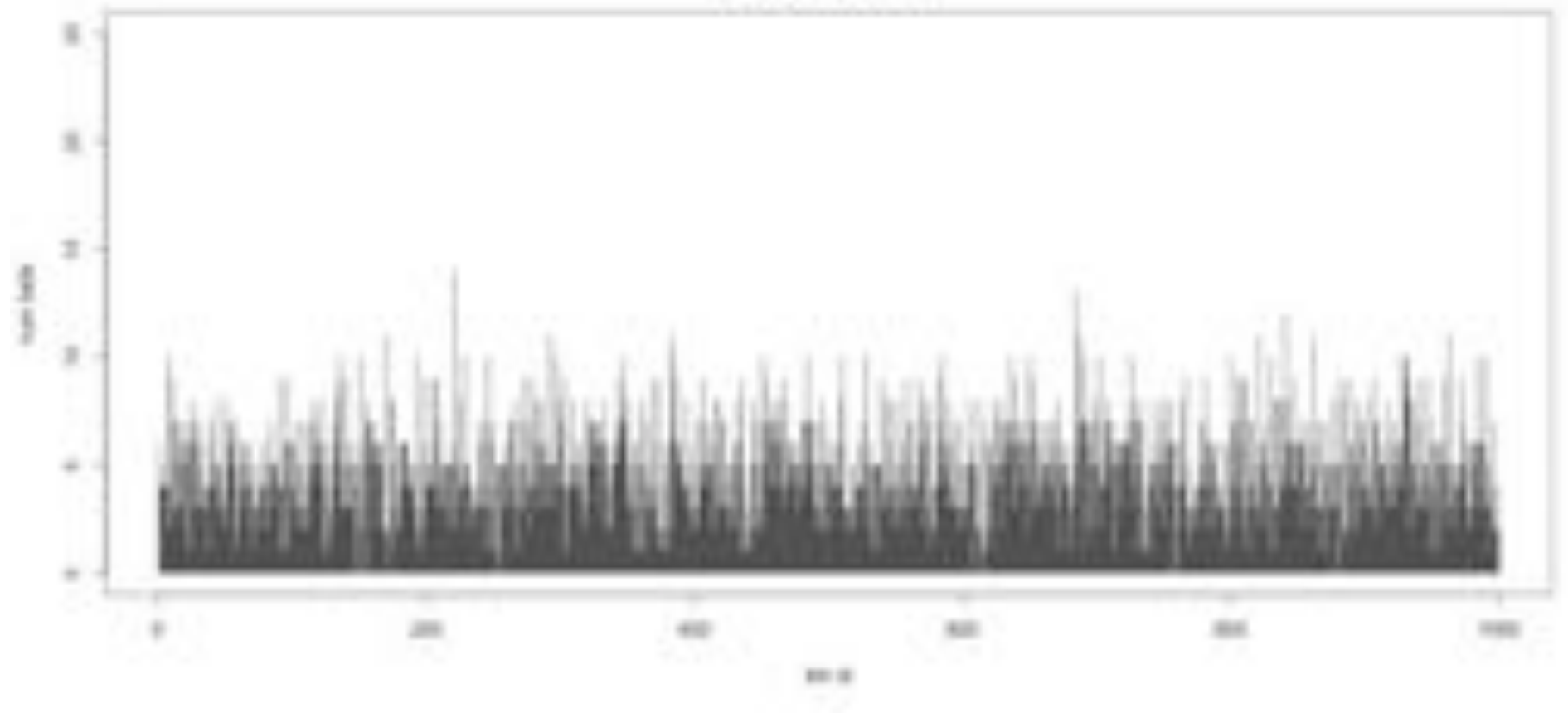


5x Sequencing

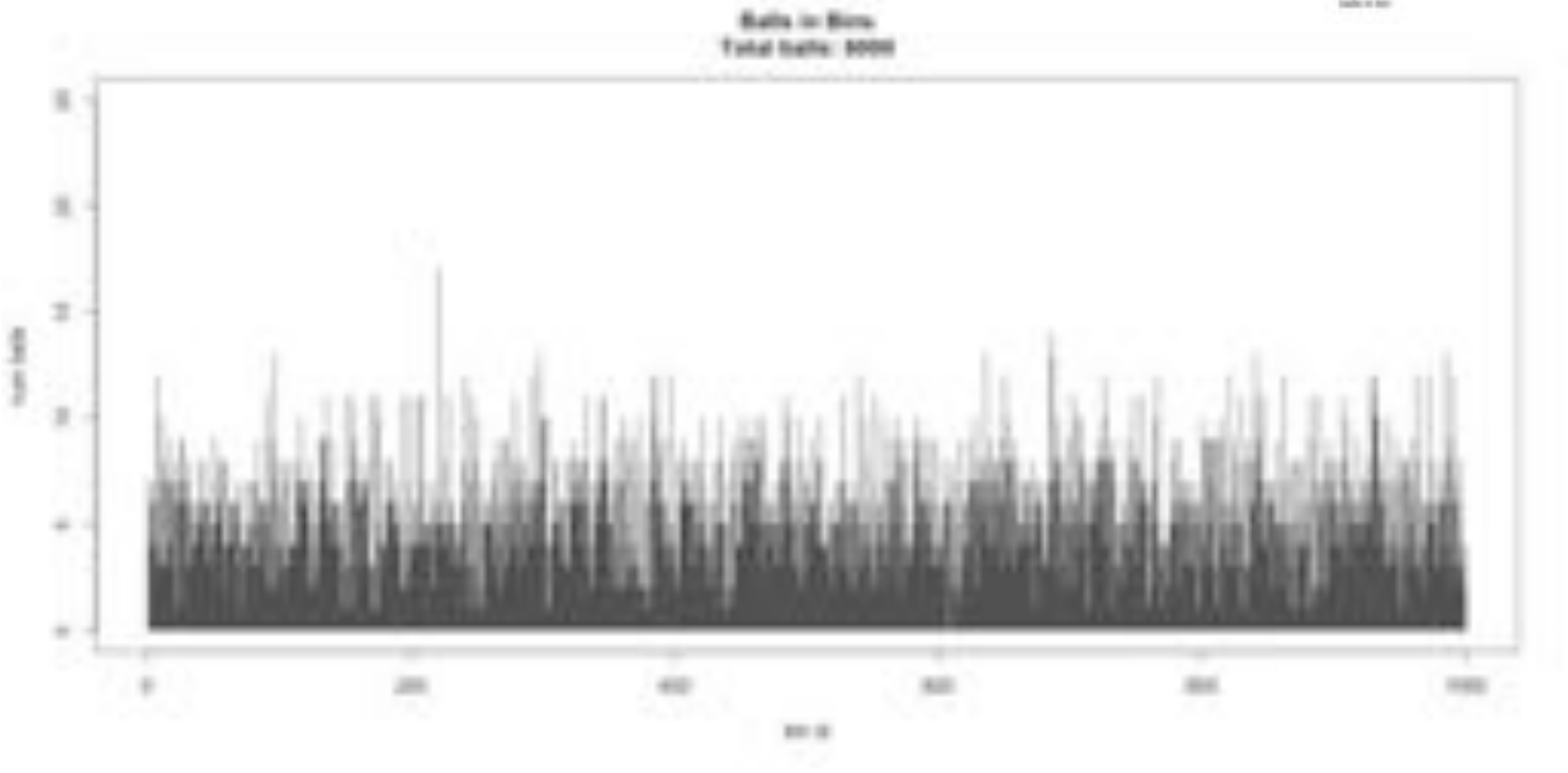
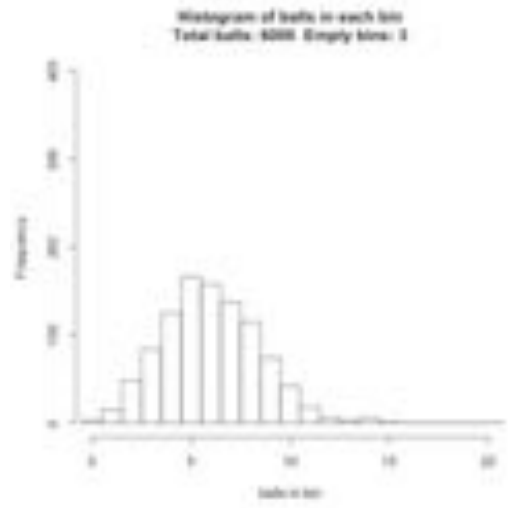
Histogram of balls in each bin
Total balls: 5000 Empty bins: 7



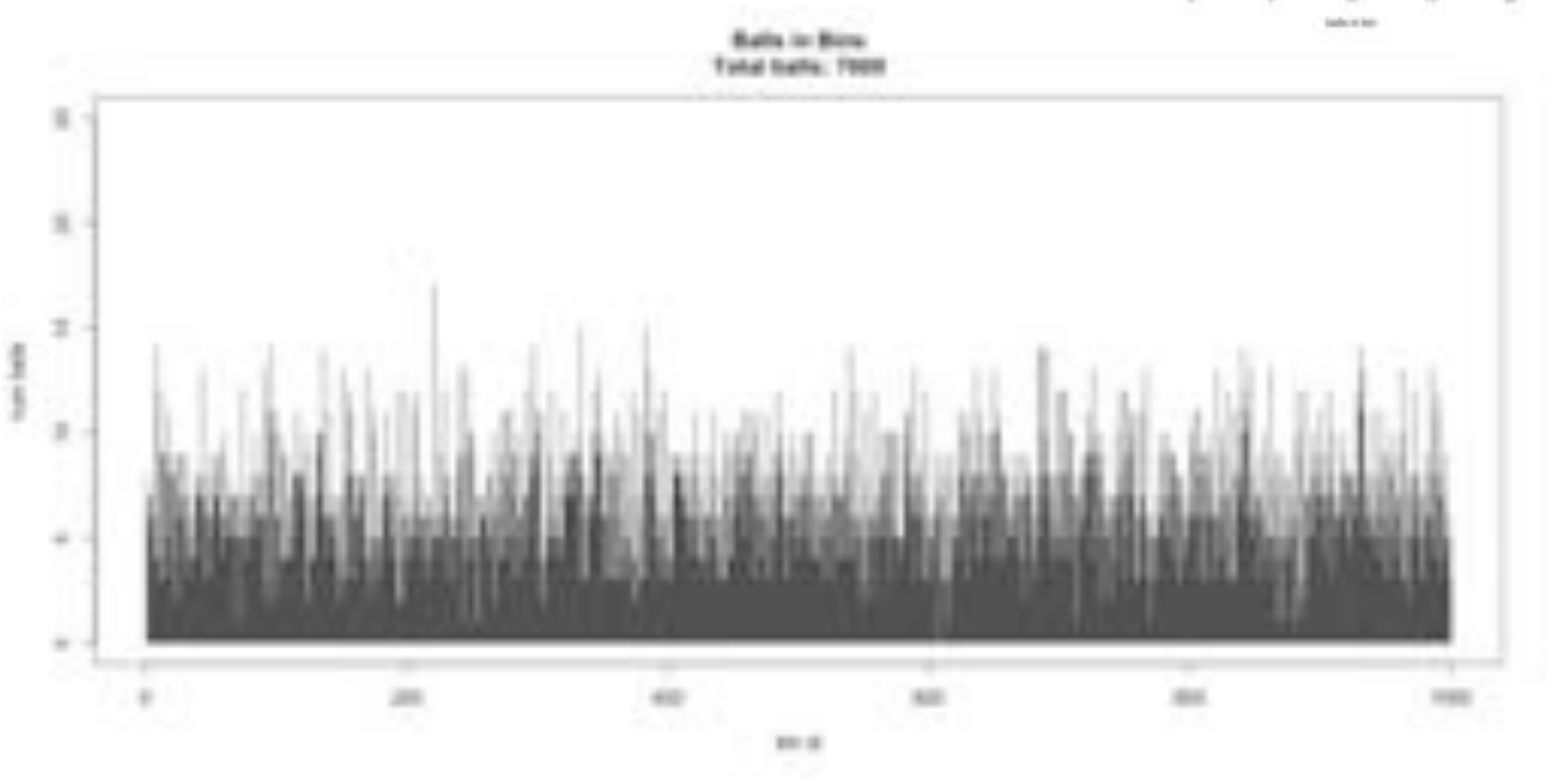
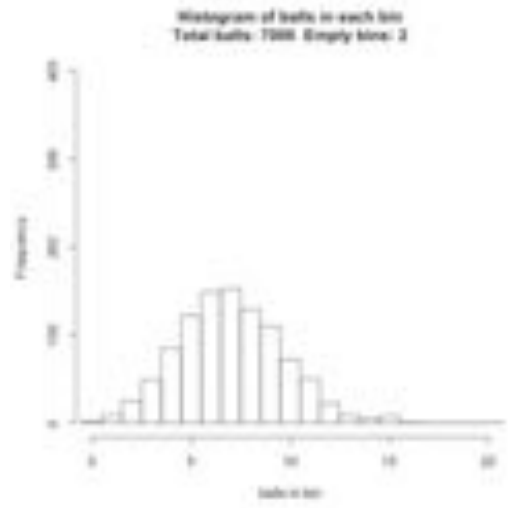
Balls in Bins
Total balls: 5000



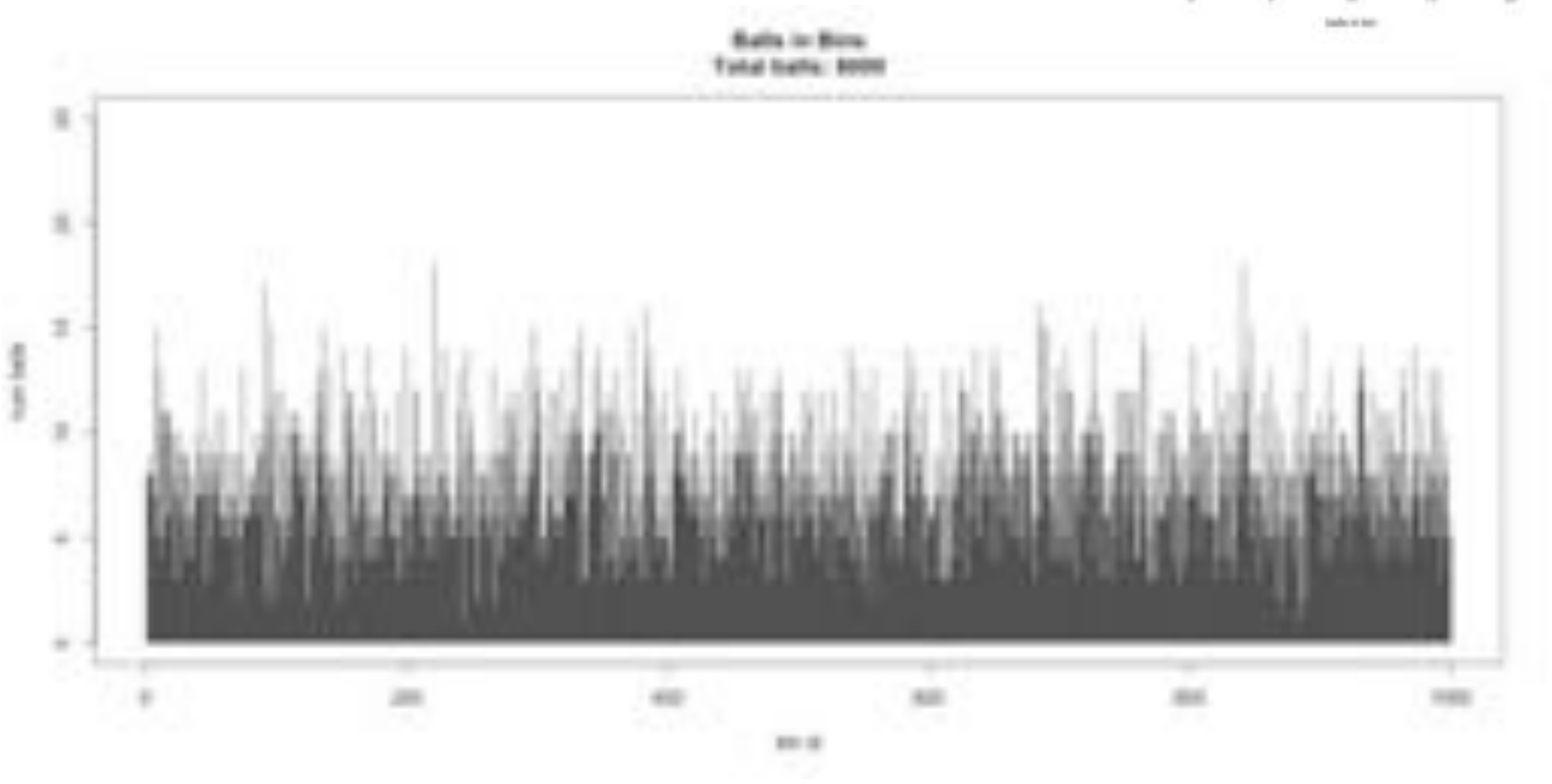
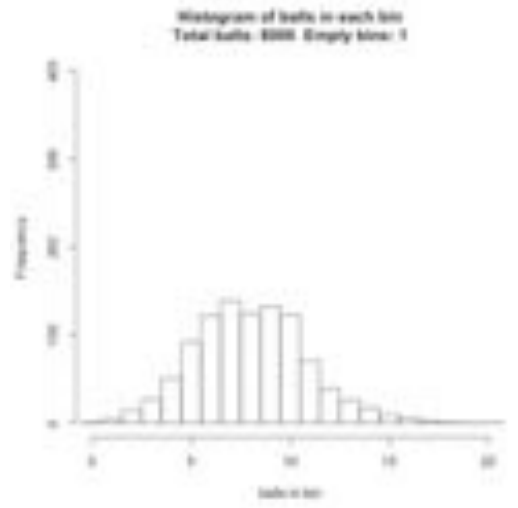
6x Sequencing



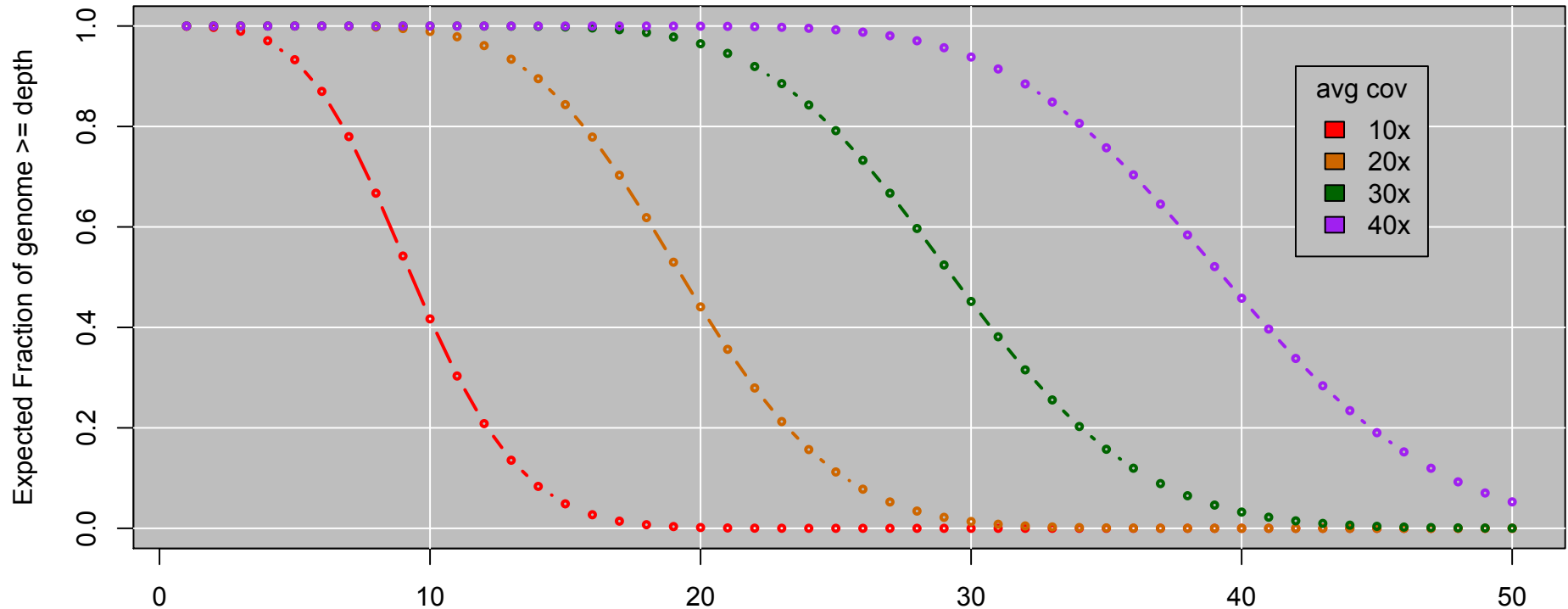
7x Sequencing



8x Sequencing



Genome Coverage Distribution

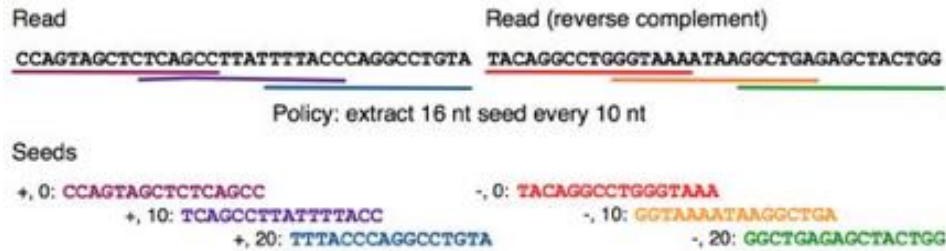


Expect Poisson distribution on depth
Standard Deviation = $\sqrt{\text{cov}}$

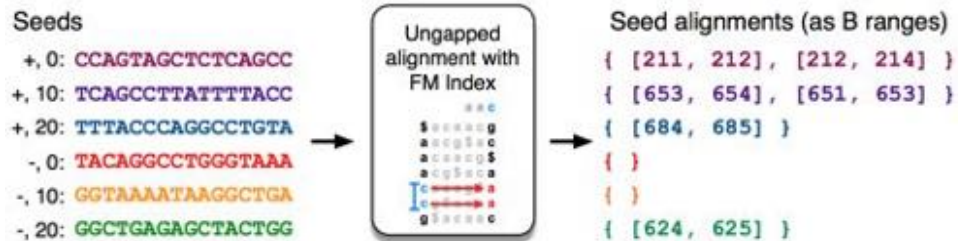
This is the mathematically model => reality may be much worse
Double your coverage for diploid genomes

Bowtie2 Overview

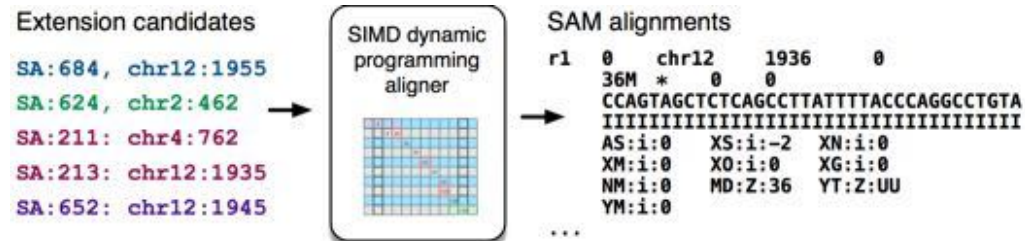
1. Split read into segments



2. Lookup each segment and prioritize



3. Evaluate end-to-end match

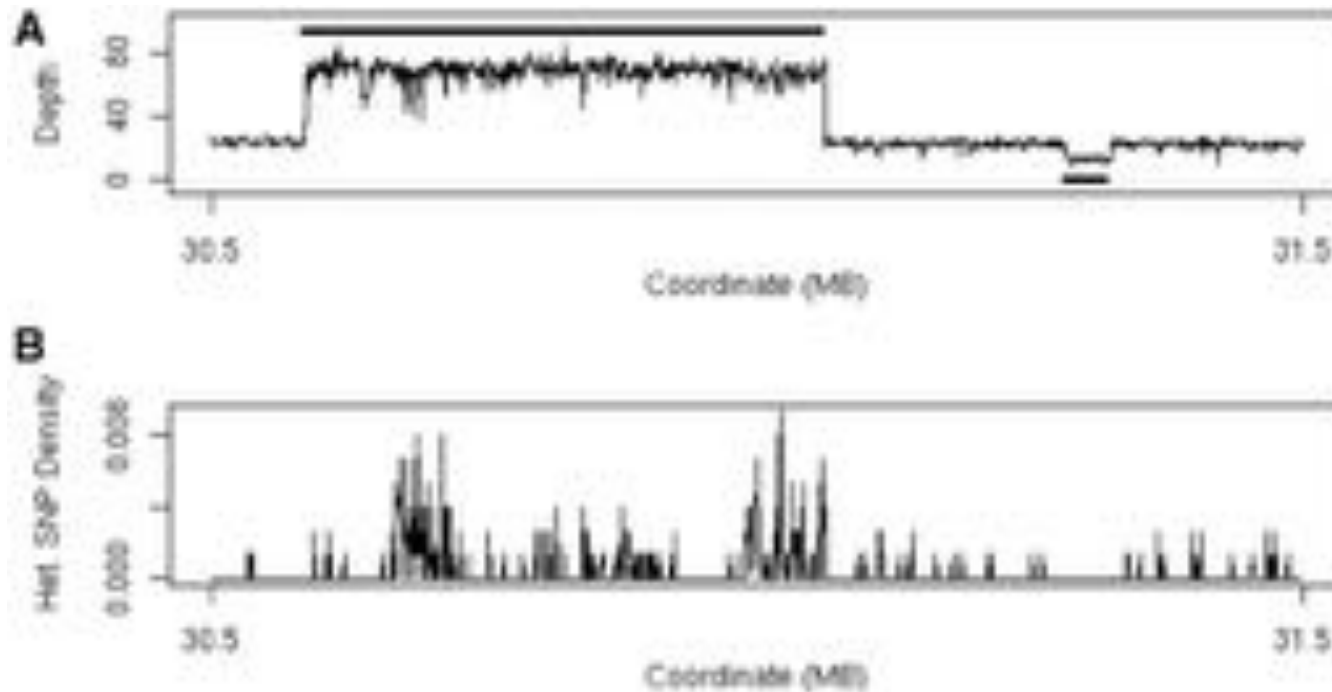


Fast gapped-read alignment with Bowtie 2.

Langmead B, Salzberg S. Nature Methods. 2012, 9:357-359.

CNV calling

Beware of (Systematic) Errors



(A) Plot of sequencing depth across a one megabase region of A/J chromosome 17 clearly shows both a region of 3-fold increased copy number (30.6–31.1 Mb) and a region of decreased copy number (at 31.3 Mb).

Simpson J T et al. *Bioinformatics* 2010;26:565-567

- Identify CNVs through increased depth of coverage & increased heterozygosity
 - Segment coverage levels into discrete steps
 - Be careful of GC biases and mapping biases of repeats

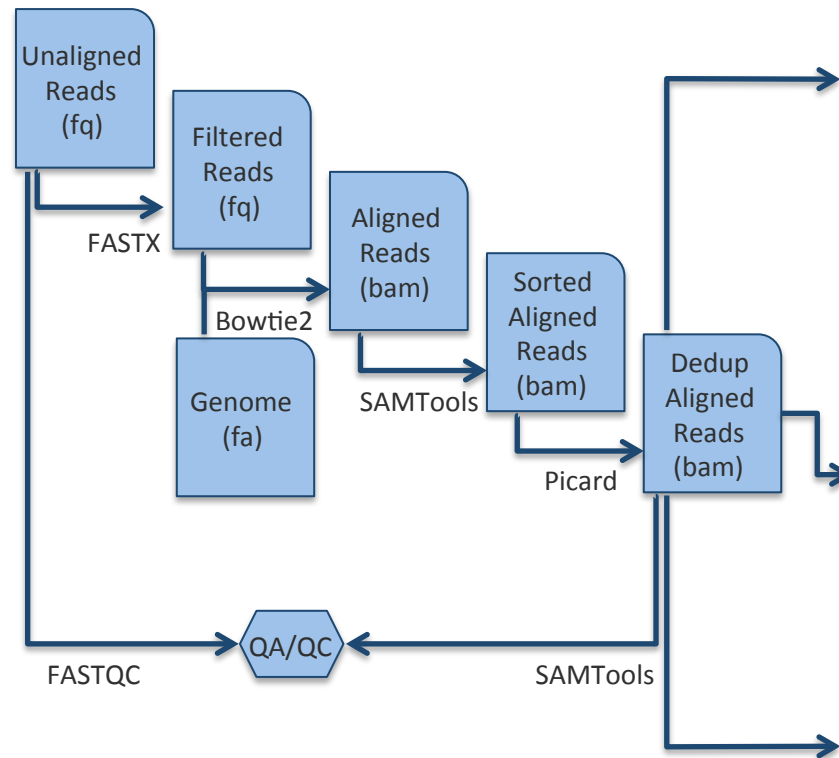
1. Introduction to KBase
2. Resequencing and variation calling theory
3. KBase services for variation calling
4. Live Demo
5. Additional Resources





Illumina HiSeq 2000
Sequencing by Synthesis

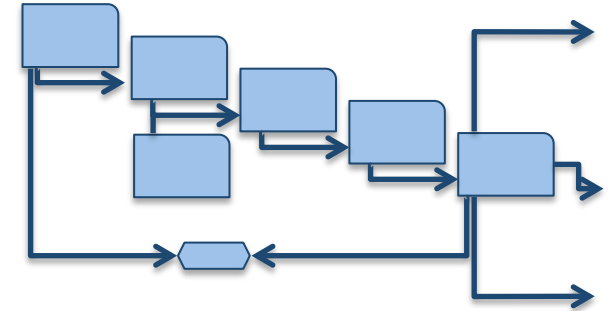
>60Gbp / day



- Assays**
- Read QA/QC
 - Mapping Stats
 - SNVs / Indels
 - CNVs / SVs
 - RNA-seq
 - ChIP-seq
 - DNase-seq
 - FAIRE-seq
 - Methyl-seq
 - ChIA-PET
 - Hi-C
 - ...

Genotyping API

- **Bowtie:** Launch alignment task with Bowtie
- **BWA:** Launch alignment task with BWA
- **SNPCalling:** Launch SNPcalling task with SAMTools
- **SortAlignments:** Launch task to sort by chromosome



Job API

- **ClusterStatus:** return basic status of cluster (jobs running, nodes available, etc)
- **JobStatus:** Given a JobID, returns current status
- **ListJobs:** List JobID running with a given username
- **KillJob:** Kills a given JobID

Data API

- **List:** List files in a directory
- **Fetch:** Fetch files from HDFS
- **Put:** Put files into HDFS
- **RM:** Delete files on HDFS
- **FetchBAM:** On-the-fly conversion to BAM
- **PutFastq:** Put reads into HDFS with conversion

Notes:

- All calls are authenticated with KBase username/password

1. Identify reference genome

```
$ all_entities_Genome -f scientific_name | grep -i 'Populus'
```

2. Upload Reads to KBase cloud

```
$ jk_fs_put_pe populus.1.fq.gz populus.2.fq.gz populus
```

3. Align Reads with Bowtie2

```
$ jk_compute_bowtie -in=populus.pe -org=populus -out=populus_align
```

4. Call SNPs with SAMTools

```
$ jk_compute_samtools_snp -in=populus_align -org=populus -out=populus_snps
```

5. Merge and Download VCF files

```
$ jk_compute_vcf_merge -in=populus_snps --alignments=populus_align -out=populus.vcf  
$ jk_fs_get populus.vcf
```

```
$ all_entities_Genome -f scientific_name | grep -i 'populus'
kb|g.3907      Populus trichocarpa
```

```
$ all_entities_Genome -f scientific_name | grep -i 'saccharomyces'
kb|g.10010     Schizosaccharomyces octosporus yfa200-2
kb|g.10042     Zygosaccharomyces bisporus IFO 1330
kb|g.10017     Schizosaccharomyces japonicus
kb|g.21735     Zygosaccharomyces rouxii
kb|g.10016     Schizosaccharomyces pombe
kb|g.2101      Saccharomyces cerevisiae S288c
kb|g.1000      Saccharomyces cerevisiae (baker's yeast)
kb|g.20490     Saccharomyces cerevisiae virus L-A (L1)
kb|g.9030      Saccharomyces cerevisiae virus L-BC (L4)
kb|g.10010     Schizosaccharomyces octosporus
kb|g.21020     Saccharomyces castellii
kb|g.20015     Saccharomyces 275 RNA narnavirus
kb|g.9739      Schizosaccharomyces japonicus yf275
kb|g.9110      Schizosaccharomyces pombe 972b-2
kb|g.10040     Zygosaccharomyces bailii
kb|g.10040     Saccharomyces cerevisiae
kb|g.1700      Schizosaccharomyces pombe
kb|g.9050      Saccharomyces servazii
kb|g.8705      Saccharomyces 205 RNA narnavirus
kb|g.21062     Saccharomyces cerevisiae r601-1a 1
kb|g.10013     Saccharomyces pastorianus Weihenstephan 34/70
kb|g.8353      Zygosaccharomyces bailii virus 2
kb|g.9401      Saccharomyces cerevisiae killer virus #1
```

Select the proper KBase ID

Identify reference genome

```
$ all_entities_Genome -f scientific_name | grep -i 'Populus'
```


Alignments



Samtools Variant Detection

SAMTools

SAMTools

SAMTools

```
##FASTQ_HEADER_17795 - T A 140
refine_B011301 13449 - C G 82
1312.0.100.01
refine_B011301 20300 - C T,G 139
20.124.01
refine_B011301 22454 - T G 88
1219.0.100.01
refine_B011301 22798 - C G 73
1208.0.100.01
refine_B011301 24089 - T A 79
1208.0.100.01
refine_B011301 25139 - G T 56
106.0.100.01
refine_B011301 26922 - T G T,8
127.0.100.01
refine_B011301 30428 - F A 71
101.0.100.01
refine_B011301 32203 - A T 141
110.0.100.01
refine_B011301 33229 - G C 76
1208.0.100.01
refine_B011301 34061 - C G 77
101.1007.0.100.01
refine_B011301 35149 - T A 55
105.0.114.01
refine_B011301 35796 - G C 55
101.100.0.104.01
refine_B011301 35798 - T TA 123
refine_B011301 35757 - T TA 108
refine_B011301 35861 - C T 82
```

```
##FASTQ_HEADER_17795 - T A 140
refine_B011301 13449 - C G 82
1312.0.100.01
refine_B011301 20300 - C T,G 139
20.124.01
refine_B011301 22454 - T G 88
1219.0.100.01
refine_B011301 22798 - C G 73
1208.0.100.01
refine_B011301 24089 - T A 79
1208.0.100.01
refine_B011301 25139 - G T 56
106.0.100.01
refine_B011301 26922 - T G T,8
127.0.100.01
refine_B011301 30428 - F A 71
101.0.100.01
refine_B011301 32203 - A T 141
110.0.100.01
refine_B011301 33229 - G C 76
1208.0.100.01
refine_B011301 34061 - C G 77
101.1007.0.100.01
refine_B011301 35149 - T A 55
105.0.114.01
refine_B011301 35796 - G C 55
101.100.0.104.01
refine_B011301 35798 - T TA 123
refine_B011301 35757 - T TA 108
refine_B011301 35861 - C T 82
```

```
##FASTQ_HEADER_17795 - T A 140
refine_B011301 13449 - C G 82
1312.0.100.01
refine_B011301 20300 - C T,G 139
20.124.01
refine_B011301 22454 - T G 88
1219.0.100.01
refine_B011301 22798 - C G 73
1208.0.100.01
refine_B011301 24089 - T A 79
1208.0.100.01
refine_B011301 25139 - G T 56
106.0.100.01
refine_B011301 26922 - T G T,8
127.0.100.01
refine_B011301 30428 - F A 71
101.0.100.01
refine_B011301 32203 - A T 141
110.0.100.01
refine_B011301 33229 - G C 76
1208.0.100.01
refine_B011301 34061 - C G 77
101.1007.0.100.01
refine_B011301 35149 - T A 55
105.0.114.01
refine_B011301 35796 - G C 55
101.100.0.104.01
refine_B011301 35798 - T TA 123
refine_B011301 35757 - T TA 108
refine_B011301 35861 - C T 82
```

Called Variants (VCF)

Call SNPs with SAMTools

```
$ jk_compute_samtools_snp -in=populus_align -org='kb|g.3907' -out=populus_snps
```


1. Identify reference genome

```
$ all_entities_Genome -f scientific_name | grep -i 'Populus'
```

2. Upload Reads to KBase cloud

```
$ jk_fs_put_pe populus.1.fq.gz populus.2.fq.gz populus
```

3. Align Reads with Bowtie2

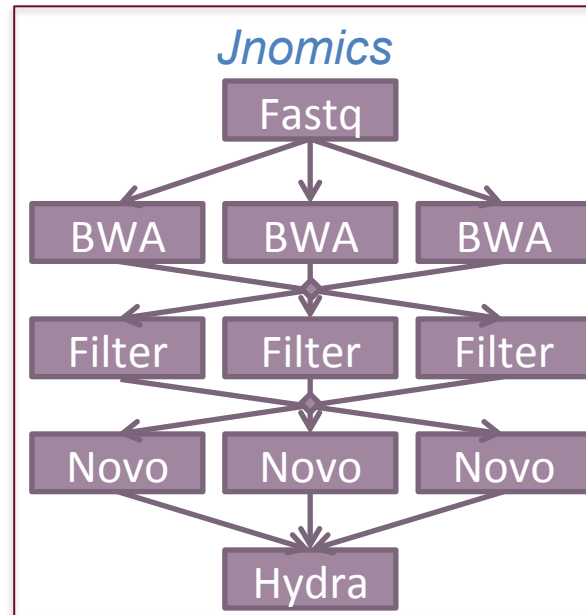
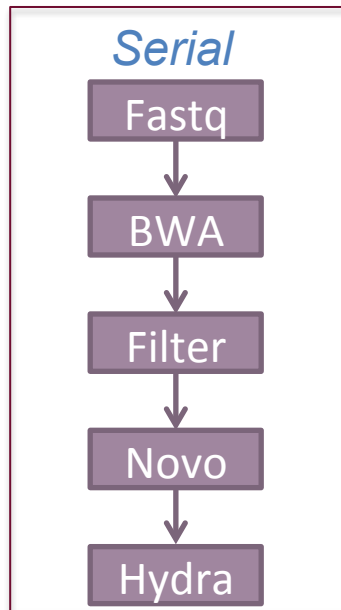
```
$ jk_compute_bowtie -in=populus.pe -org=populus -out=populus_align
```

4. Call SNPs with SAMTools

```
$ jk_compute_samtools_snp -in=populus_align -org=populus -out=populus_snps
```

5. Merge and Download VCF files

```
$ jk_compute_vcf_merge -in=populus_snps --alignments=populus_align -out=populus.vcf  
$ jk_fs_get populus.vcf
```

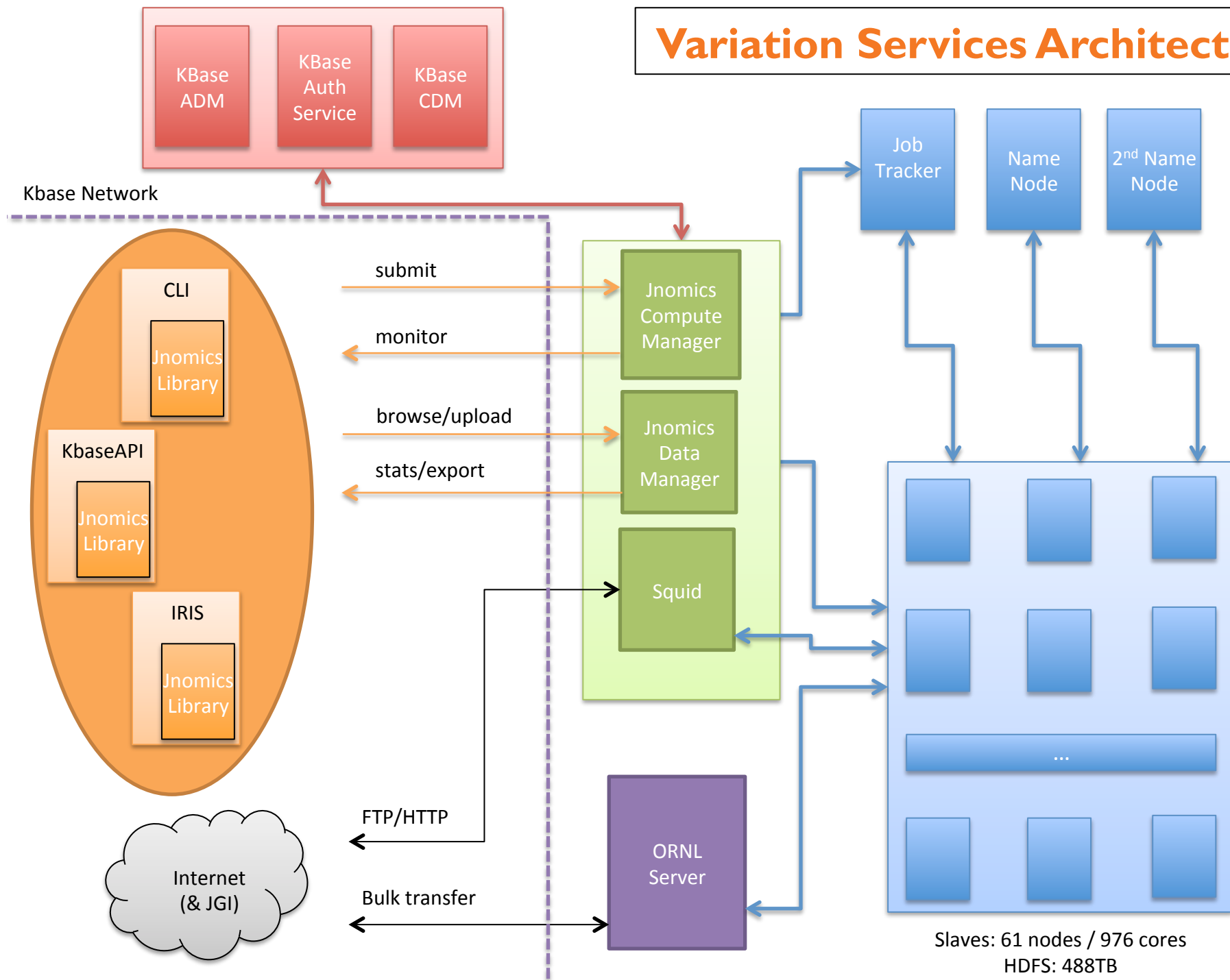


- Rapid parallel execution of data-intensive analysis
 - FASTX, BWA, Bowtie2, Novoalign, SAMTools, Hydra
 - Sorting, merging, filtering, selection, clustering, correlating
 - Supports BAM, SAM, BED, fastq

Answering the demands of digital genomics

Titmus, MA, Gurtowski, J, Schatz, MC (2012) *Concurrency & Computation*

Variation Services Architecture



Align & call SNPs from 35M 80bp (14Gbp) reads with maize genome (zmb73v2)
 Identified 372k high confidence SNPs

	Serial	Multicore	KBase Cloud
Config	1 core (1 node)	44 core (1 node)	118 cores (15 nodes)
Bowtie2	45 h*	1h 10m	23 m
Sort	2 hr	2 hr	N/A
Samtools	2 hr	2 hr	12 m
End-to-End	50h*	5h 10m	35 m
Speedup	1x	9.6x	86x

*estimated time

Maize Population Analysis

Align & call SNPs from 131 maize samples
1 TB fastq / 408Gbp input data

	Serial	KBase cloud (small)	KBase Cloud (large)
Config	1 core (1 node)	210 cores (15 nodes)	854 cores (61 nodes)
Bowtie2	1311 hr*	19.5 hr	5 hr
Sort	58 hr*	N/A	N/A
Samtools	58 hr*	3.5 hr	1.5 hr
End-to-End	1427 hr*	23 hr	6.5 hr
Speedup	1x	62x	219x

*estimated time

1. Introduction to KBase
2. Resequencing and variation calling theory
3. KBase services for variation calling
4. Live Demo
5. Additional Resources



Online Demo

1. Browse to KBase website: <http://kbase.us/>
2. Sign up for KBase account: <https://gologin.kbase.us/SignUp>
3. Download KBase DMG: <http://kbase.us/for-users/get-started/>
Or use IRIS: <http://kbase.us/services/docs/invocation/Iris/>
4. Variation Services Tutorial:
<http://kbase.us/for-users/tutorials/analyzing-data/variation-service/>
5. Summarize mutations:

```
$ cat yeast.vcf  
$ grep -v '^#' yeast.vcf | cut -f1 | sort | uniq -c  
$ grep -v '^#' yeast.vcf | cut -f 4,5 | sort | uniq -c | sort -nrk1 | head
```

1. Introduction to KBase
2. Resequencing and variation calling theory
3. KBase services for variation calling
4. Live Demo
5. Additional Resources



Resource	URL
KBase	http://kbase.us/
Getting Started	http://kbase.us/for-users/user-home/
Variation Services	http://kbase.us/for-users/tutorials/analyzing-data/variation-service/
Bowtie2	http://bowtie-bio.sourceforge.net/bowtie2/index.shtml
BWA	http://bio-bwa.sourceforge.net/
SAMTools	http://samtools.sourceforge.net/
VCF Spec	http://www.1000genomes.org/wiki/Analysis/Variant%20Call%20Format/vcf-variant-call-format-version-40
SNPeff	http://snpeff.sourceforge.net/
KBase Contact	http://kbase.us/contact-us/
Survey	https://www.surveymonkey.com/s/KB-user-info

Thank You!

<http://schatzlab.cshl.edu>
@mike_schatz / @DOEKBase

