# Genome Skimming Workshop:
## Using Low Coverage WGS to Estimate Distances and Sequencing Parameters.

by Eduardo Charvel and Siavash Mirarab

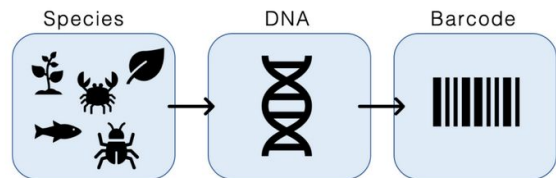# Cloning GitHub and Installing Tools:

```
git clone https://github.com/echarvel3/skimming_scripts-echarvel.git
```



```
cd skimming_scripts-echarvel
bash ./install.sh
```

# Intro: Evolutionary Ecology Approaches

## Genetic Markers



Species — DNA — Barcode

**marker genes, mtDNA, microsatellites**

**PROS:**
- cost-effective (single genes)
- no recombination (mtDNA)

**CONS:**
- expensive (larger bait sets)
- limited resolution (single genes)
- reference dependency
- homoplasy
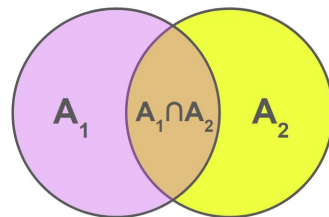
## SNP based

**RAD-seq, ANGSD, GATK**

**PROS:**
- genome-wide resolution (WGS)

**CONS:**
- reference-based (WGS)
- higher coverage ( >10x ) (WGS)
- restriction-enzymes limit resolution (RAD-seq)



## *k*-mer based



$A_1$   $A_1 \cap A_2$   $A_2$
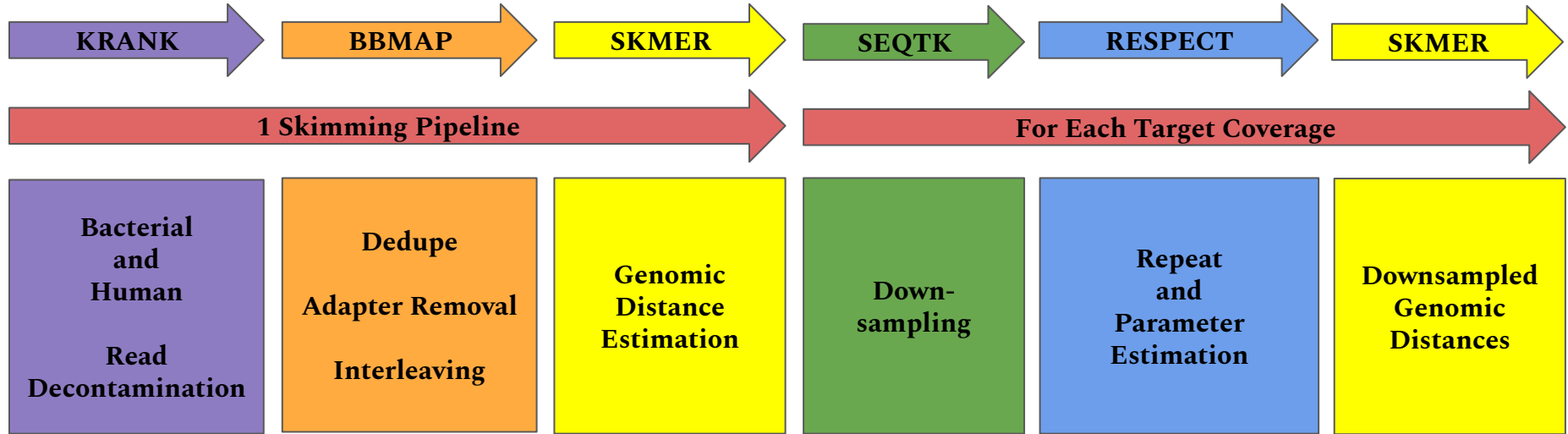
**PROS:**
- genome-wide resolution
- low coverage ( < 1x .. 8x )
- alignment-free

**CONS:**
- sensitive to sequencing quality:
  - contamination
    (*addressed in this workshop*)
  - repetitiveness
    (*active area of research*)
  - hybridization

# Skimming Pipeline Overview

| KRANK → | BBMAP → | SKMER → | SEQTK → | RESPECT → | SKMER → |
|---|---|---|---|---|---|

| 1 Skimming Pipeline → | | | For Each Target Coverage → | | |
|---|---|---|---|---|---|
| Bacterial and Human<br><br>Read Decontamination | Dedupe<br><br>Adapter Removal<br><br>Interleaving | Genomic Distance Estimation | Down-sampling | Repeat and Parameter Estimation | Downsampled Genomic Distances |

- **Meant to easily perform all operations necessary for accurate genomic distance and sequencing parameter estimation.**

# BBtools preprocessing:



*merged reads*

*interleaved reads*

| 1 | **ATG**CTACC...T |
|---|---|
| 2 | **ATG**GGAA...A |
| 3 | **ATG**TTTA...C |
| 4 | **ATG**AAAA...A |
| 5 | **ATG**CTGC...C |
| 6 | **ATG**GGCA...G |
| 7 | **ATG**GACC...T |
| 8 | **ATG**GACC...T |
| 9 | **ATG**GACC...T |
| 10 | **ATG**GACC...T |

| 1 | CTACC...T |
|---|---|
| 2 | GGAA...A |
| 3 | TTTA...C |
| 4 | AAAA...A |
| 5 | CTGC...C |
| 6 | GGCA...G |
| 7 | GACC...T |
| 8 | GACC...T |
| 9 | GACC...T |
| 10 | GACC...T |

Sequencing Adapters (bbduk.sh)

Deduplicate (dedupe.sh)

join read1 and read2

5

# Running BBMAP

```
cd skimming_scripts-echarvel
conda activate ${CONDA_ENV} #replace with env name
cd ./test/
gunzip ./skims/*gz
mkdir bbmap_reads
bash ../bbmap_pipeline.sh \
        ./skims/read1.fq \
        ./skims/read2.fq \
        ./bbmap_reads/read_out.fq
```

## Main Skmer Paper

# Skmer: assembly-free and alignment-free sample identification using genome skims

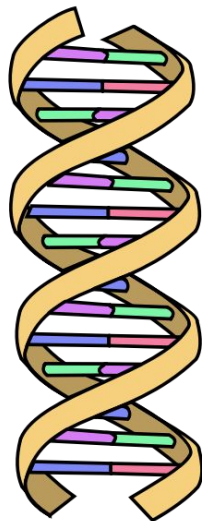Shahab Sarmashghi, Kristine Bohmann, M. Thomas P. Gilbert, Vineet Bafna ✉ & Siavash Mirarab ✉

## Skmer Bootstrapping Approach

# Quantifying the uncertainty of assembly-free genome-wide distance estimates and phylogenetic relationships using subsampling

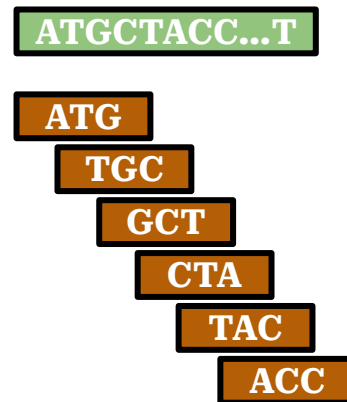Eleonora Rachtman [1], Shahab Sarmashghi [2], Vineet Bafna [3], Siavash Mirarab [2 4]

7

# SKMER Theory:

collection of
reads ~ 150 bp

each can be decomposed
into a set of *k*-mers

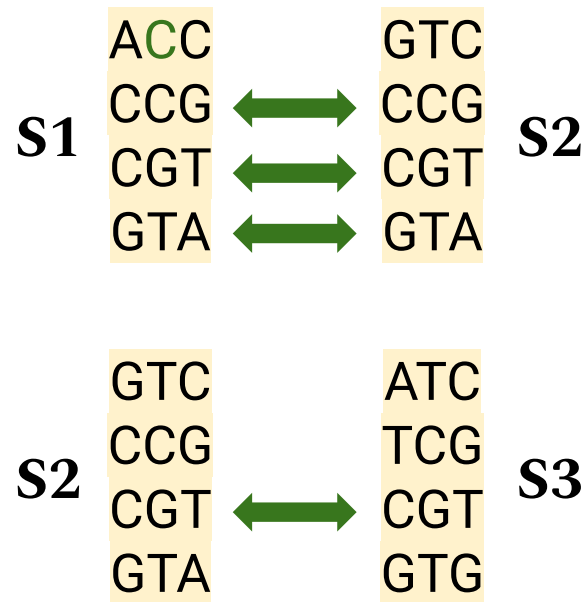| | |
|---|---|
| 1 | ATGCTACC...T |
| 2 | CTAGGAA...A |
| 3 | CCCTTTA...C |
| 4 | GATAAAA...A |
| 5 | TATACTGC...C |
| 6 | GTCGGCA...G |
| 7 | GGACTGC...C |
| 8 | CTTATCC...G |
| 9 | TAAGTGG...C |
| 10 | CAGGACC...T |

ATGCTACC...T

ATG
TGC
GCT
CTA
TAC
ACC

Obtain Biological Samples ⟩ Library Preparation ⟩ Short Read Sequencing ⟩ *k*-mer counting

# SKMER Theory:

# Mash (2016): $k$-mer Based Distances



**shared $k$-mers = intersection =** $|A_1 \cap A_2|$

**set of all $k$-mers = union =** $|A_1 \cup A_2|$

**Jaccard Similarity Index =** $J = \dfrac{|A_1 \cap A_2|}{|A_1 \cup A_2|}$

$$D = 1 - \left(\frac{2J}{1+J}\right)^{1/k}$$
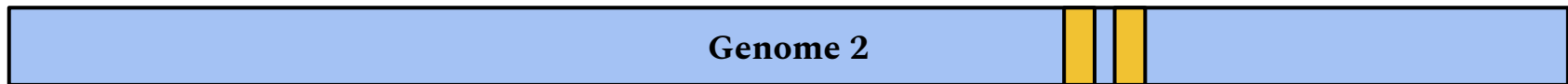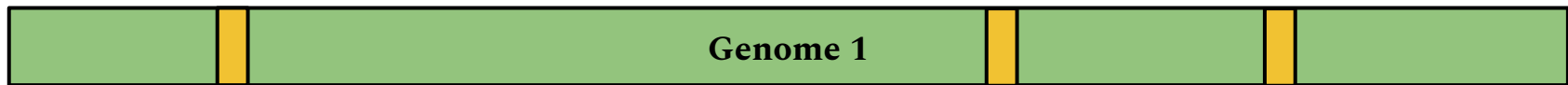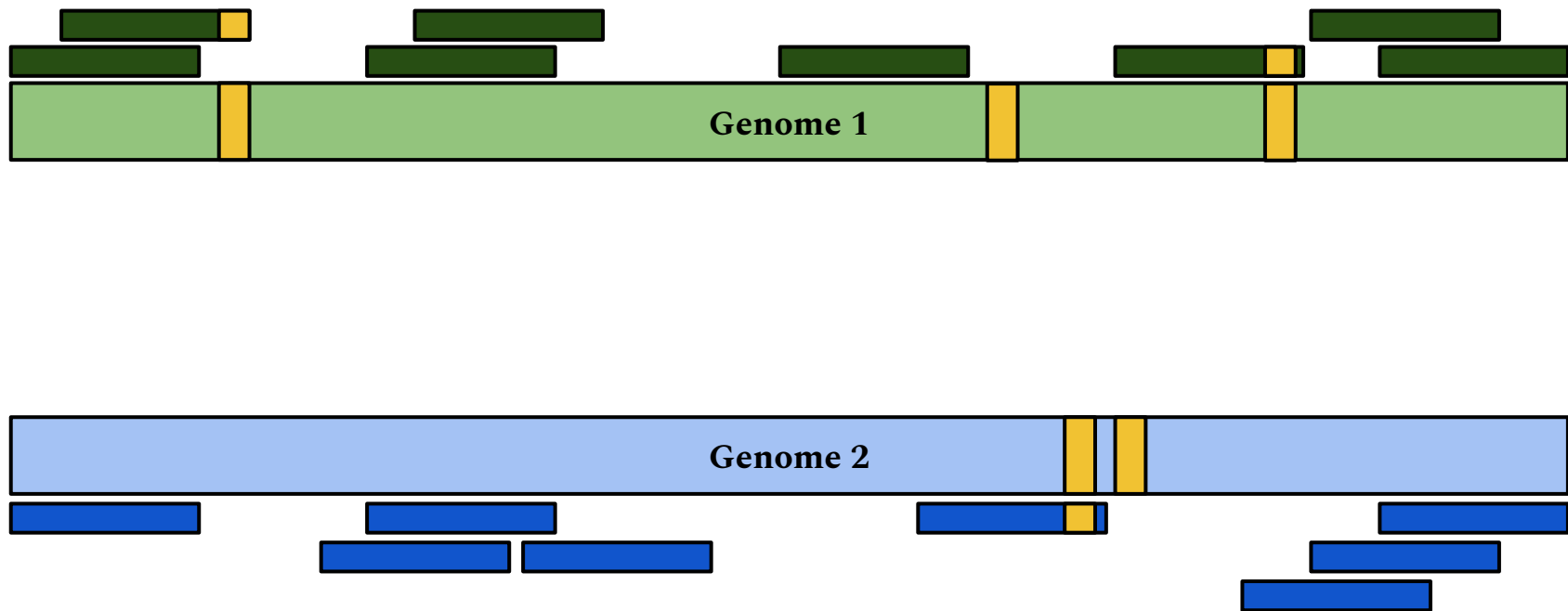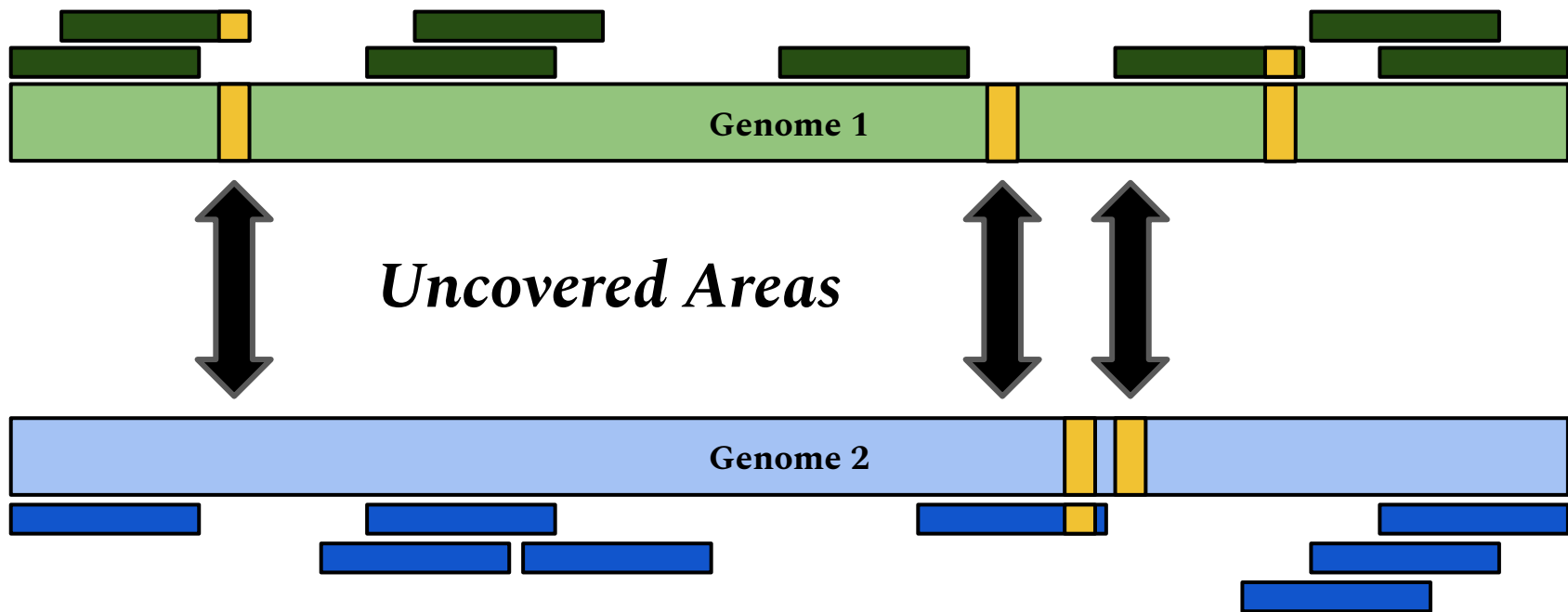
# SKMER: Using Low Coverage Skims

Genome 1

Genome 2

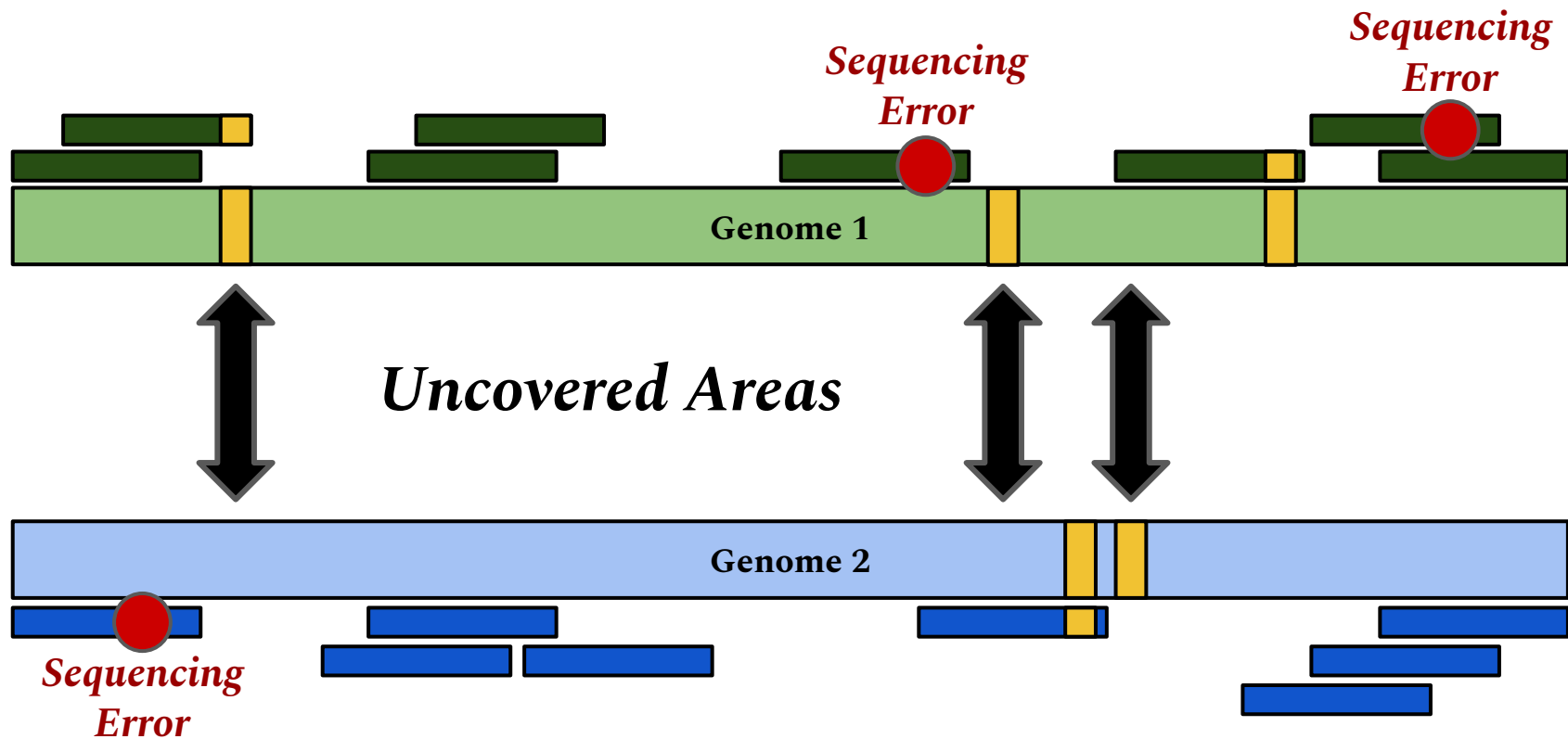# SKMER: Using Low Coverage Skims

# SKMER: Using Low Coverage Skims

# SKMER: Using Low Coverage Skims

# SKMER: Using Low Coverage Skims

# SKMER: Using Low Coverage Skims

$$D = 1 - \left(\frac{2J}{1+J}\right)^{1/k}$$

$$D = 1 - \left(\frac{2(\zeta_1 L_1 + \zeta_2 L_2)J}{\eta_1 \eta_2 (L_1 + L_2)(1+J)}\right)^{1/k}$$

$\eta_i$ (eta) = probability a $k$-mer is covered at least once without error.

$\zeta_i$ (zeta) = the total number of k-mers observed from both genomes.

$L_i$ = genome size.

Both **eta** and **zeta** are functions of **coverage** and **error**.

# SKMER: Estimating Parameters

**λ** (lambda) = how many times a *k*-mer is covered

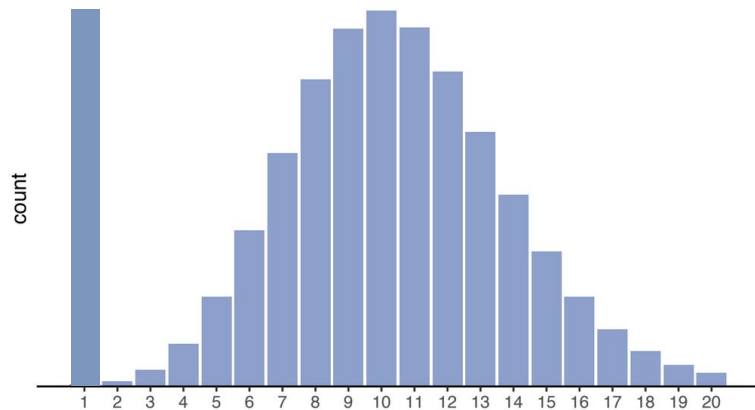**ξ** (xi) =  number of error free reads covering a *k*-mer

**ε** (epsilon) = sequencing error rate

**|A∩B|** (obs. intersection) = shared *k*-mers between samples.

- Use *k*-mer frequency spectrum to estimate of parameters.

- Use *MASH* to compute observed intersection.

# Estimating Error and Coverage:



$$P(X) = \frac{\lambda^x e^{-\lambda}}{X!}$$

WGS follows a Poisson Distribution

$h = \mathrm{argmax}_{i \geq 2} M_i$ = *mode*

$\xi = \frac{M_{h+1}}{M_h}(h+1)$ = *number of error-free reads covering a k-mer*

$\lambda = \frac{M_1}{M_h}\frac{\xi^h}{h!}e^{-\xi} + \xi(1 - e^{-\xi})$ = *k-mer coverage*

$\epsilon = 1 - (\xi/\lambda)^{1/k}$ = *sequencing error rate*

# Running Skmer Reference

```
cd skimming_scripts-echarvel
conda activate ${CONDA_ENV} #replace with env name
gunzip ./test/skims/*gz
skmer reference ./skims/bbmap_reads -p1
```

# Other Skmer Utilities:

```
skmer 3.2.1 - Estimating genomic distances between genome-skims

optional arguments:
  -h, --help            show this help message and exit
  --debug               Print the traceback when an exception is raised

commands:
  reference    Process a library of reference genome-skims or assemblies
  distance     Compute pairwise distances for a processed library
  query        Compare a genome-skim or assembly against a reference library
  subsample    Performs  subsample on a library of reference genome-skims or assemblies
  correct      Performs correction of subsampled distance matrices obtained for reference genome-skims or assemblies

  {reference,subsample,correct,distance,query}
                        Run skmer {commands} [-h] for additional help
```
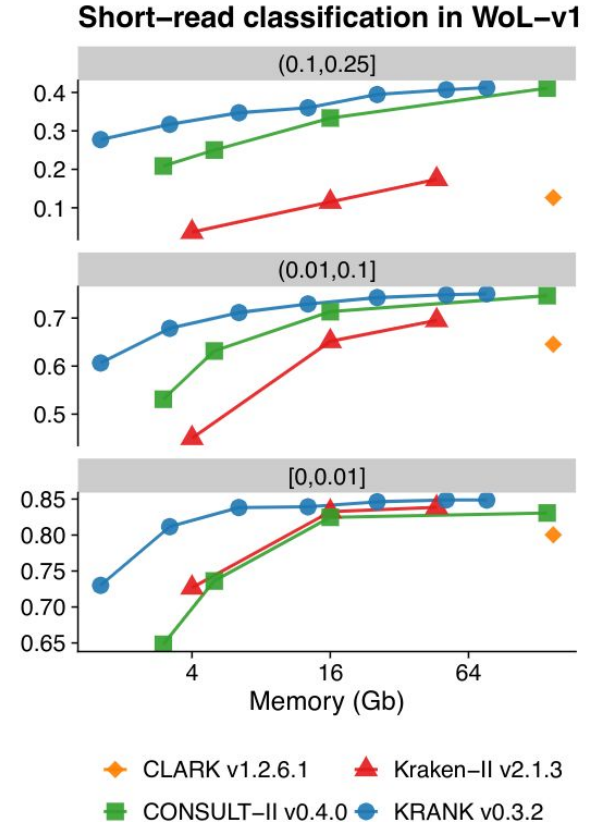
*__distance__ == generate distance matrix (automatic for reference)*
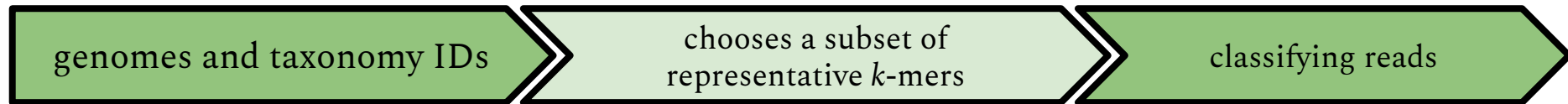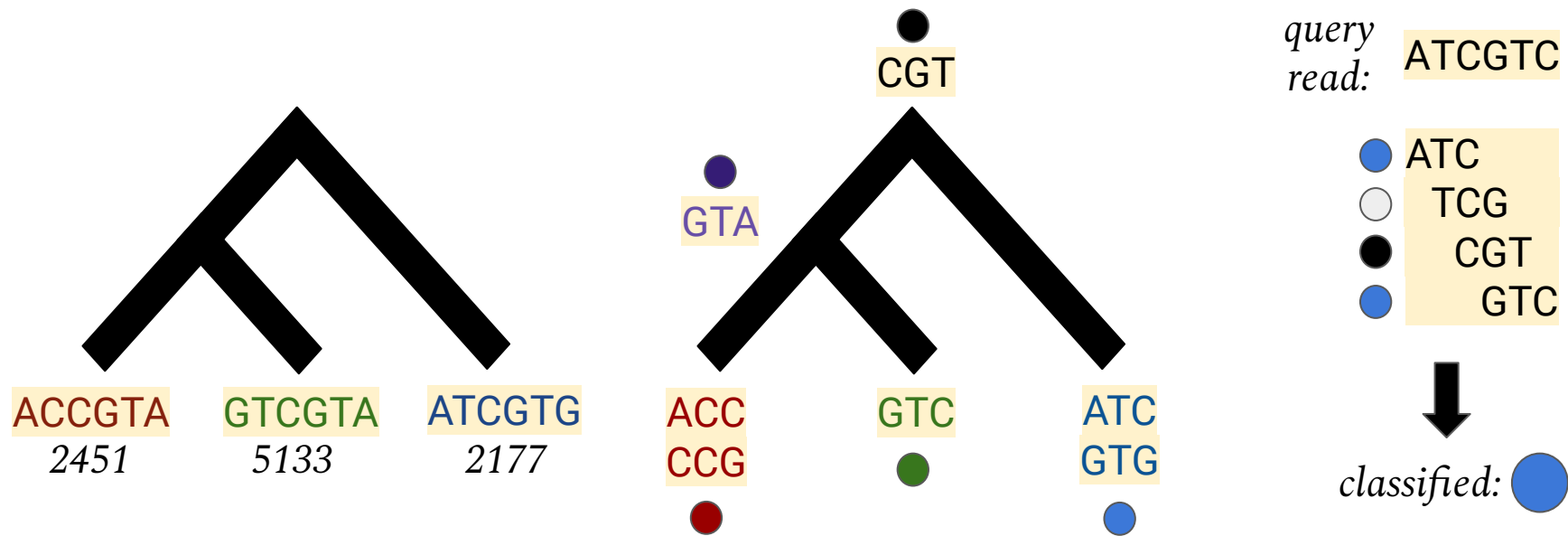*__query__ == add new skim/assembly to library*
*__subsample__ == perform bootstrap approach*

# Memory-bound *k*-mer selection for large and evolutionary diverse reference libraries

Ali Osman Berk Şapcı, Siavash Mirarab

**Short–read classification in WoL–v1**

Legend:
- CLARK v1.2.6.1
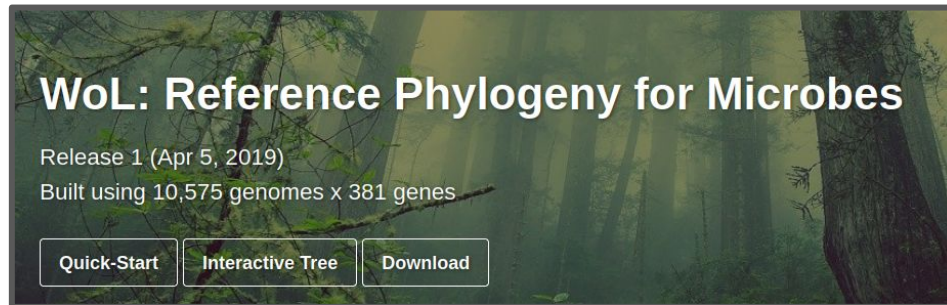- Kraken–II v2.1.3
- CONSULT–II v0.4.0
- KRANK v0.3.2

# KRANK: Decontamination

# KRANK: Decontamination

Most Common Sources of Contamination:

- **Bacterial** → WoL-v1 dataset - archaeal and bacterial genomes



- **Human** → T2T Pangenome [05-2024]

# KRANK: Decontamination

**genomes**

**tax IDs**

| | | | |
|---|---|---|---|
| 1 | ATGCTACC...T | | 52631 |
| 2 | CTAGGAA...A | | 23100 |
| 3 | CCCTTTA...C | | 23100 |
| 4 | GATAAAA...A | × | 48019 |
| 5 | TATACTGC...C | | 5001 |
| 6 | GTCGGCA...G | | 662 |
| 7 | GGACTGC...C | | 80919 |
| 8 | CTTATCC...G | | 80919 |

```
krank build
    -l $LIBRARY_DIRECTORY
    -t $TAXONOMY_DIRECTORY
    -i $MAPPING_FILE
    --from-library
    --batch-size 8
    --target-batch 0
    --num-threads $NUM_THREADS
```

**https://github.com/bo1929/KRANK**

genomes and taxonomy IDs

pass genomes and tax IDs to KRANK

# Running KRANK build

```
cd ./skimming_scripts-echarvel/test/KRANK_test
```

krank build \
  -l ${LIBDIR} -t ./taxonomy/ \
  -i ./input_map.tsv \
  -k 27 -w 35 -h 12 -b 8 -s 2\
  --from-scratch --input-sequences \
  --kmer-ranking representative --adaptive-size --lca soft \
  --num-threads ${NTHREADS} \


krank build \
  -l ${LIBDIR} -t ./taxonomy/ \
  -i ./input_map.tsv \
  -k 27 -w 35 -h 12 -b 8 -s 2\
  --target-batch 0 --fast-mode --from-library --input-sequences --keep-intermediate \
  --kmer-ranking representative --adaptive-size --lca soft \
  --num-threads ${NTHREADS}

# Running KRANK query

```
cd ./skimming_scripts-echarvel/test/
```

```
${SCRIPT_DIR}/KRANK/krank query \
        --library-dir ${LIBRARIES} \
        --query-file ./bbmap_read.fq  \
        --max-match-distance 5 \
        --total-vote-threshold 0.03 \
        --num-threads ${NUM_THREADS} \
        --output-dir "${OUTPUT_DIRECTORY}/krank_output/krank_reports/"
```
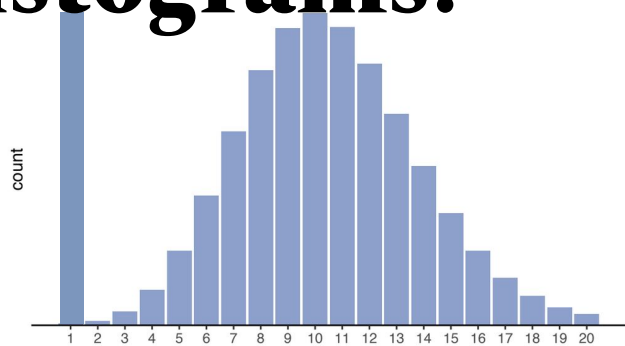
# Estimating repeat spectra and genome length from low-coverage genome skims with RESPECT

## coverage genome skims with RESPECT

Shahab Sarmashghi, Metin Balaban, Eleonora Rachtman, Behrouz Touri, Siavash Mirarab, Vineet Bafna ✉
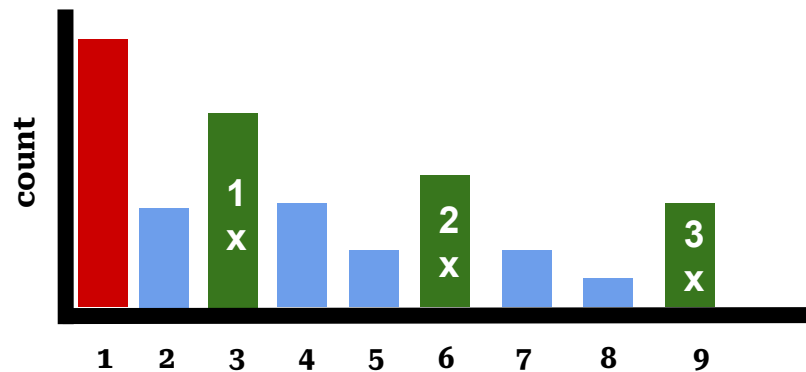
# How Repeats Affect $k$-mer Histograms:
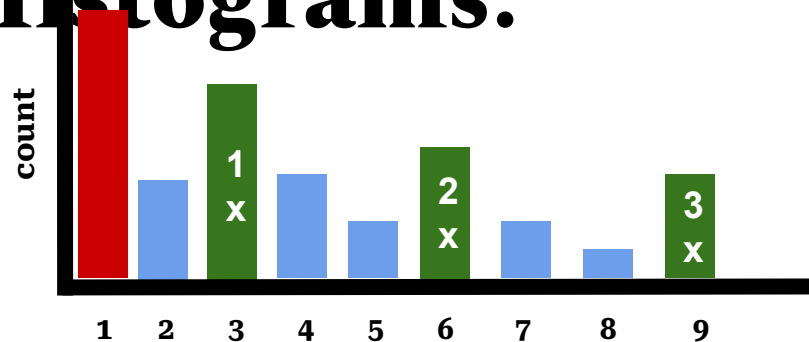


simple genome skim



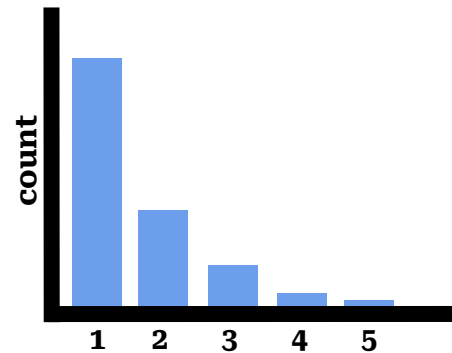repetitive genome skim

Multiple peaks and coverages.

**Can we get measures of repetitiveness using this?**

# How Repeats Affect *k*-mer Histograms:



repetitive genome

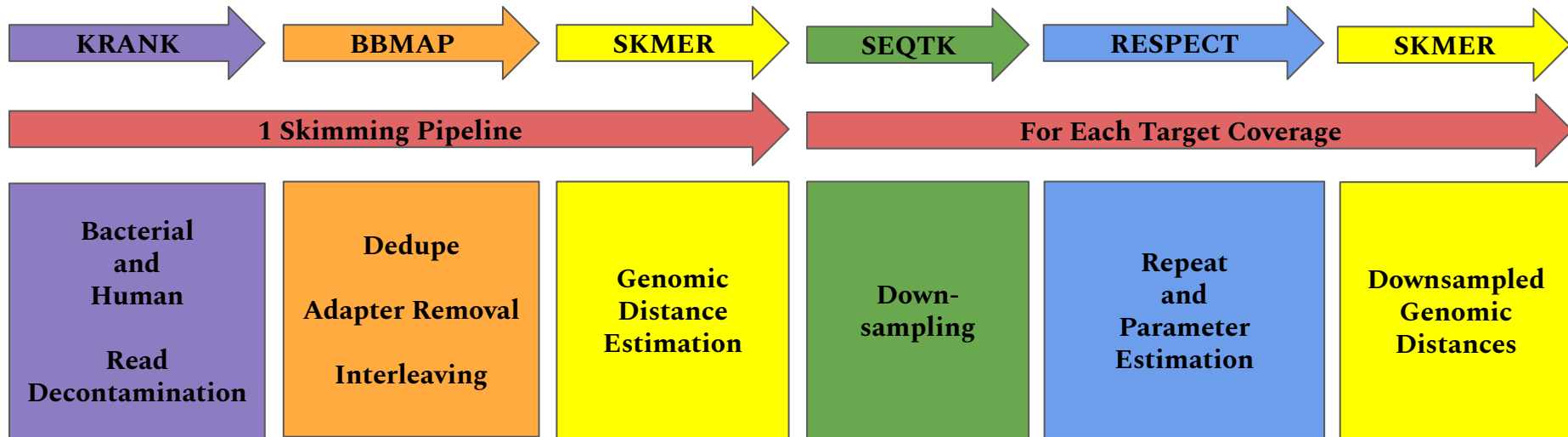Repeat spectra of an assembly

**RESPECT output:**
- repeat-informed coverage, error, and genome length estimates
- repeat spectra of a genome
- repetitiveness metrics:
    - Uniqueness Ratio, High Copy Repeats per Million

# Running RESPECT

```
cd ./skimming_scripts-echarvel/test/
respect -d ./bbmap_reads/ -N 2 --threads 1
```

# Using the Skimming Pipeline

| KRANK | BBMAP | SKMER | SEQTK | RESPECT | SKMER |
|-------|-------|-------|-------|---------|-------|

| 1 Skimming Pipeline | For Each Target Coverage |
|---------------------|--------------------------|

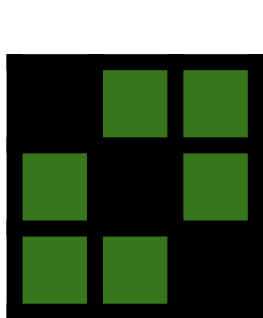| Bacterial and Human<br><br>Read Decontamination | Dedupe<br><br>Adapter Removal<br><br>Interleaving | Genomic Distance Estimation | Down-sampling | Repeat and Parameter Estimation | Downsampled Genomic Distances |
|---|---|---|---|---|---|

```
bash ./skimming_pipeline.sh -h
```

# APPLES: Scalable Distance-Based Phylogenetic Placement with or without Alignments FREE

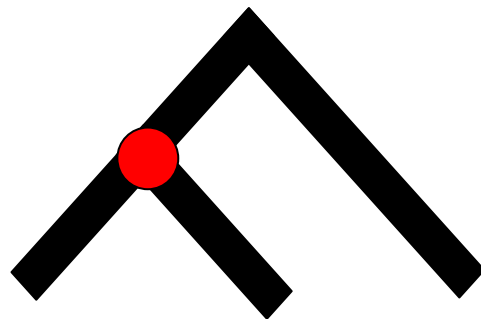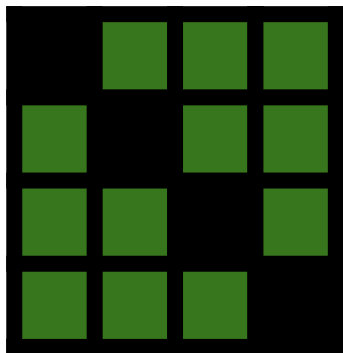Metin Balaban, Shahab Sarmashghi, Siavash Mirarab ✉

# Installing Apples:

```
### Instal APPLES
python -m pip install -U apples
run_apples.py -h
python -m pip  list |grep apples
### If you have versions older than 1.3.0, you may need to updating using:
python -m pip install --upgrade apples
```



skmer query



Apples placement

# Thank you!

**Siavash Mirarab and Vineet Bafna**

Mirarab Lab:
Isaac, Daira, Homere, Ali, Shayasteh, Nora, and Yueyu

Field Museum: Grainger Bioinformatics Center
Minderoo Foundation