# RNA virus full genome sequencing and haplotype reconstruction

Sebastian Krautwurst

February 13, 2018
33rd TBI Winterseminar in Bled
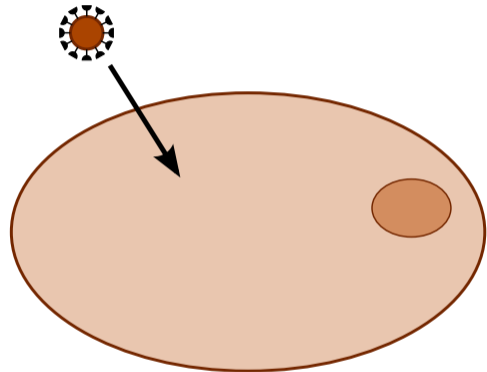
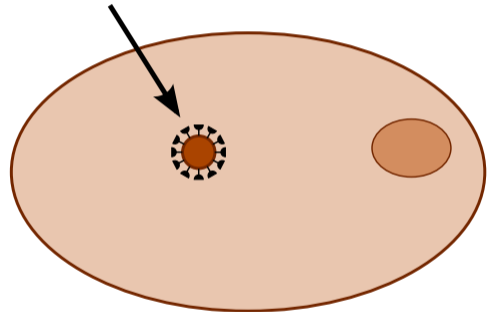FRIEDRICH-SCHILLER-
**UNIVERSITÄT
JENA**

# Background

# VIRAL HAPLOTYPES

- One species = one genome?

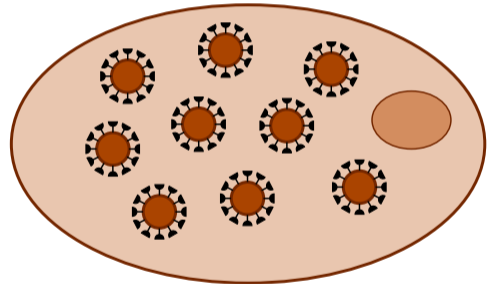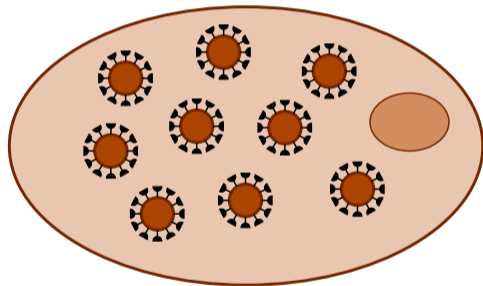# VIRAL HAPLOTYPES

- One species = one genome?

# VIRAL HAPLOTYPES

▶ One species = one genome?

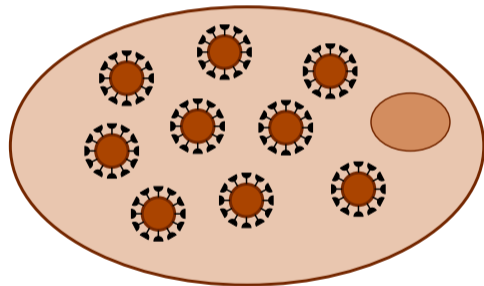# VIRAL HAPLOTYPES

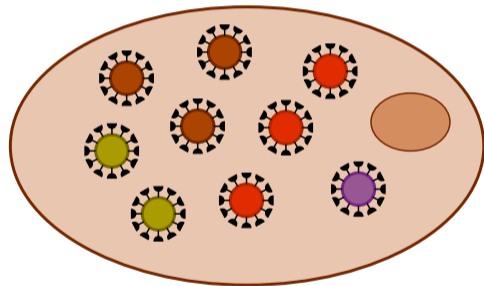- One species = one genome?
- RNA viruses: error-prone replication

# VIRAL HAPLOTYPES

- One species = one genome?
- RNA viruses: error-prone replication
- Mutation, recombination, segment reassortment

# VIRAL HAPLOTYPES

- One species = one genome?
- RNA viruses: error-prone replication
- Mutation, recombination, segment reassortment
- Diverse spectrum of genomes
  ⇒ Quasispecies reconstruction



FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

# NANOPORE SEQUENCING

- ONT MinION

# NANOPORE SEQUENCING

▶ ONT MinION



nanoporetech.com/sites/default/files/s3/minion-cutout.png

# NANOPORE SEQUENCING

- ▶ ONT MinION
- ▶ 10-20 Gb per flow cell



nanoporetech.com/sites/default/files/s3/minion-cutout.png

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

## NANOPORE SEQUENCING

- ONT MinION
- 10-20 Gb per flow cell
- Very long reads possible –
  up to 1 Mb



nanoporetech.com/sites/default/files/s3/minion-cutout.png

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

## NANOPORE SEQUENCING

- ▶ ONT MinION
- ▶ 10-20 Gb per flow cell
- ▶ Very long reads possible – up to 1 Mb
- ▶ Noisy – 15% indels



nanoporetech.com/sites/default/files/s3/minion-cutout.png
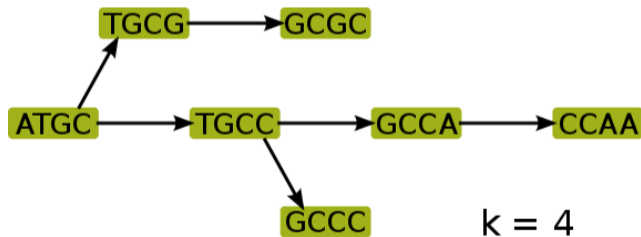
# NANOPORE SEQUENCING

- ▶ ONT MinION
- ▶ 10-20 Gb per flow cell
- ▶ Very long reads possible – up to 1 Mb
- ▶ Noisy – 15% indels
- ▶ Direct RNA sequencing protocol kit


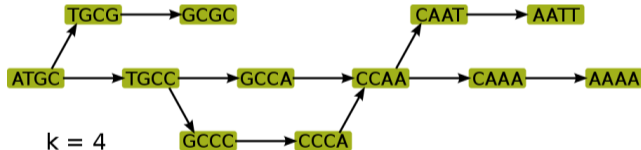
nanoporetech.com/sites/default/files/s3/minion-cutout.png

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

# DE BRUIJN GRAPH
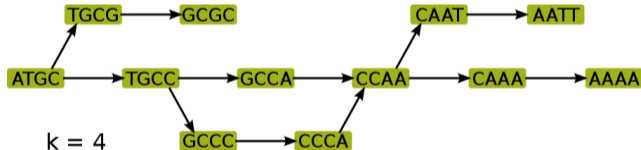
▶ Constructed from
  overlapping k-mers



k = 4

# DE BRUIJN GRAPH

- Constructed from overlapping k-mers
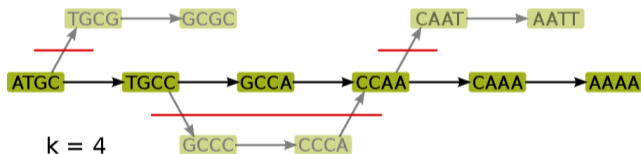- Captures variants

# DE BRUIJN GRAPH

- ▶ Constructed from overlapping k-mers
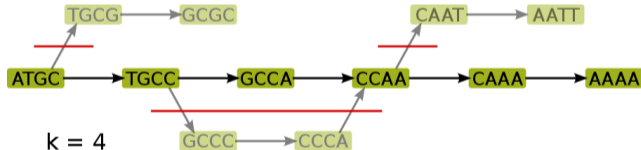- ▶ Captures variants
- ▶ Assembly: consensus



k = 4

# DE BRUIJN GRAPH

- ▶ Constructed from overlapping k-mers
- ▶ Captures variants
- ▶ Assembly: consensus
- ▶ Tip- and bulge removal



k = 4

# DE BRUIJN GRAPH

- Constructed from overlapping k-mers
- Captures variants
- Assembly: consensus
- Tip- and bulge removal
- Collapse unambiguous chains
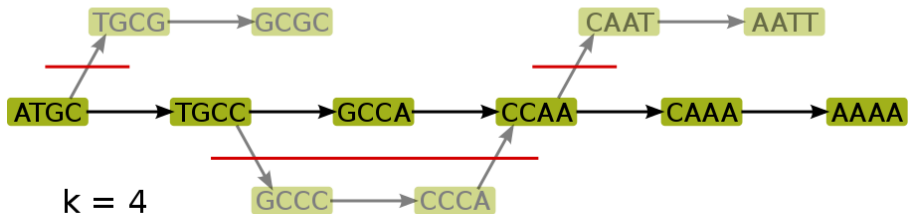


$k = 4$

# ASSEMBLY BY DE BRUIJN GRAPH

► Established de novo assembly method (Velvet, SPAdes)

# ASSEMBLY BY DE BRUIJN GRAPH

- ▶ Established de novo assembly method (Velvet, SPAdes)
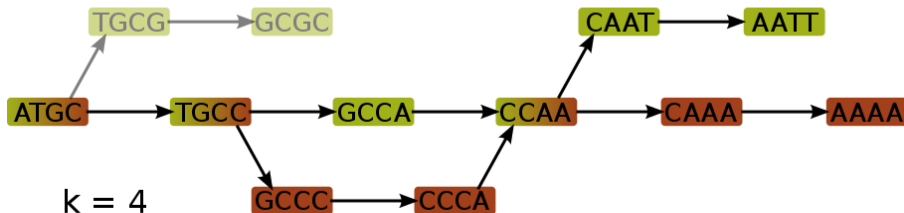- ▶ Goal: Enhance with focus on quasispecies reconstruction
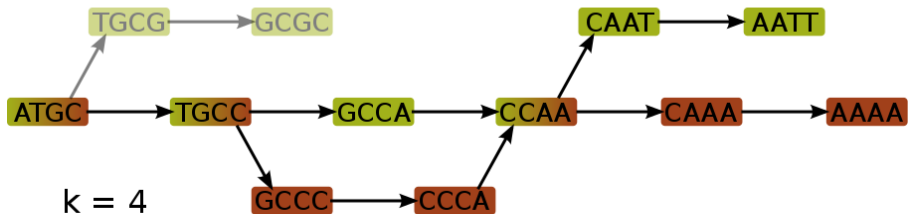


k = 4

# ASSEMBLY BY DE BRUIJN GRAPH

- ▶ Established de novo assembly method (Velvet, SPAdes)
- ▶ Goal: Enhance with focus on quasispecies reconstruction
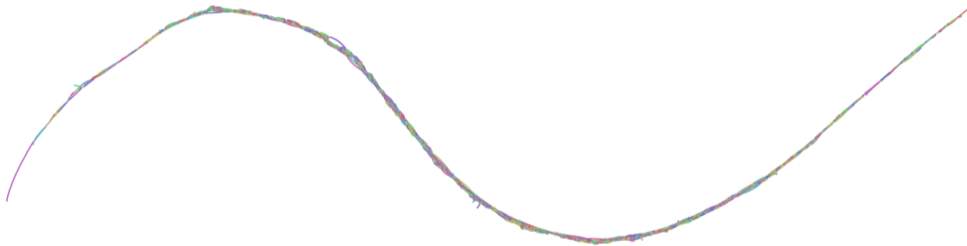- ▶ Separate haplotypes by graph manipulation, long read information
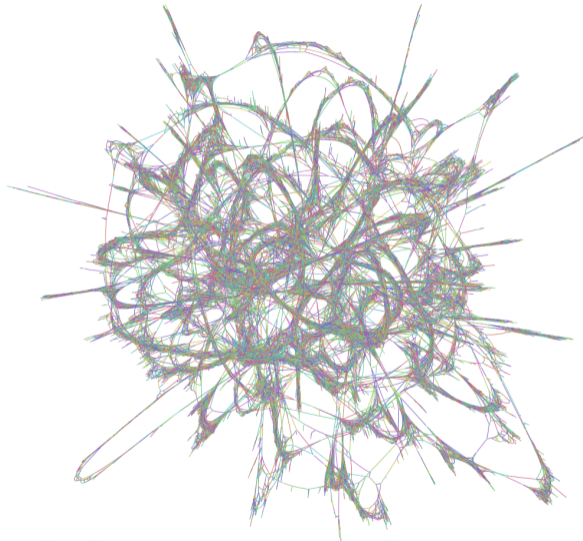
# ASSEMBLY BY DE BRUIJN GRAPH

- ▶ Established de novo assembly method (Velvet, SPAdes)
- ▶ Goal: Enhance with focus on quasispecies reconstruction
- ▶ Separate haplotypes by graph manipulation, long read information
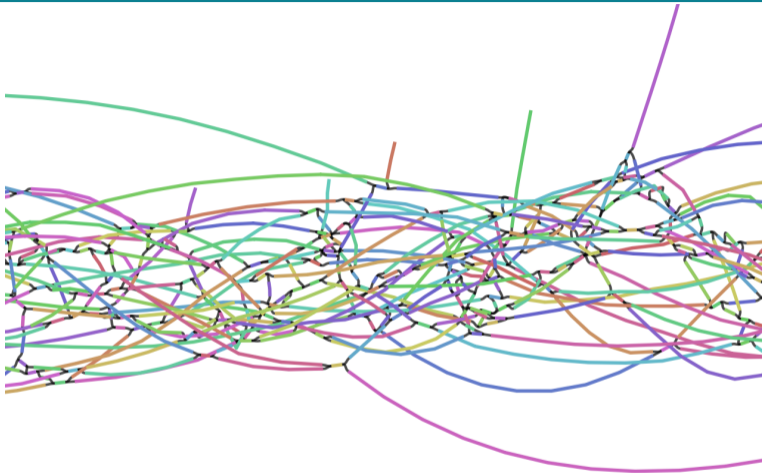- ▶ Assemble haplotype consensus sequences
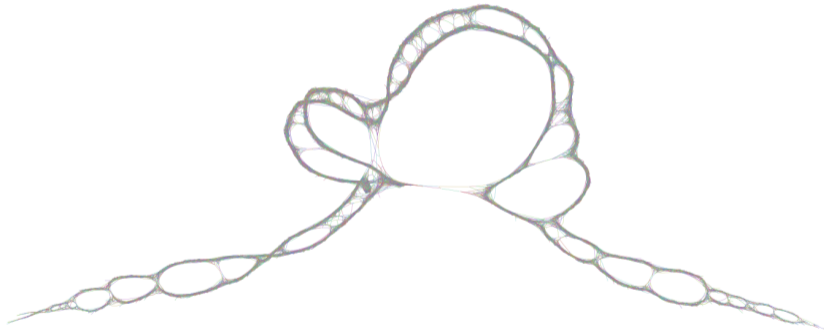
# Results so far

corona simul - 1 genome - 20 reads - 10000 nt - 10 % indels - k=25

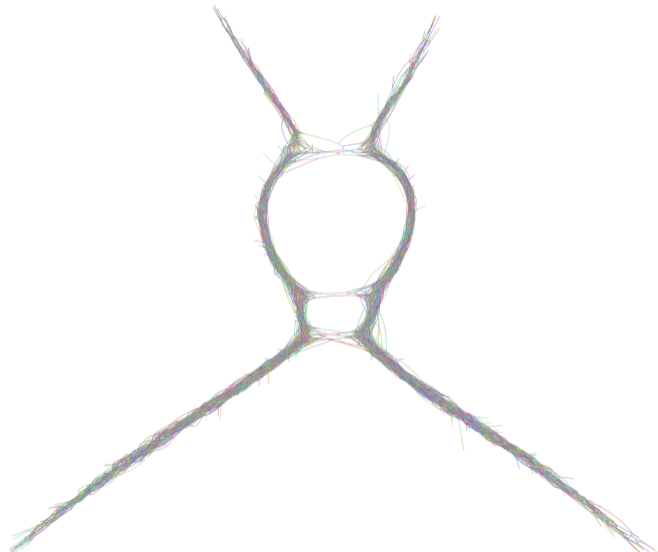FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

corona simul - 1 genome - 100 reads - 10000 nt - 10 % indels - **k=16**

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

corona - 1 genome - 40 reads - 10000 nt - **15 % indels** - k=30

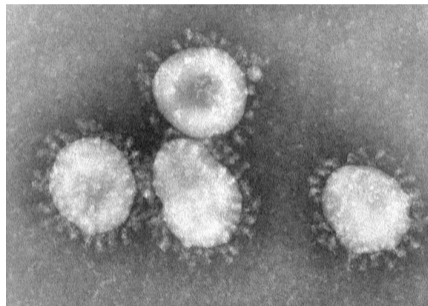FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

bvdv simul - **2 genomes** - 100/100 reads - 6000 nt - 10 % indels - k=20

bvdv simul - **2 genomes** - 100/100 reads - 6000 nt - 10 % indels - **k=30**

# REAL CORONAVIRUS READ DATA
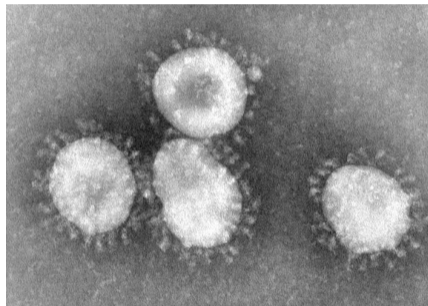
- ▶ HCoV 229E in human cell culture



en.wikipedia.org/wiki/Coronavirus#/media/

File:Coronaviruses_004_lores.jpg

# REAL CORONAVIRUS READ DATA

- HCoV 229E in human cell culture
- Direct RNA protocol kit on MinION



`en.wikipedia.org/wiki/Coronavirus#/media/`

`File:Coronaviruses_004_lores.jpg`

# REAL CORONAVIRUS READ DATA

- HCoV 229E in human cell culture
- Direct RNA protocol kit on MinION
- 293406 reads, 27 % virus, rest human



```
en.wikipedia.org/wiki/Coronavirus#/media/
```

```
File:Coronaviruses_004_lores.jpg
```

# REAL CORONAVIRUS READ DATA

- HCoV 229E in human cell culture
- Direct RNA protocol kit on MinION
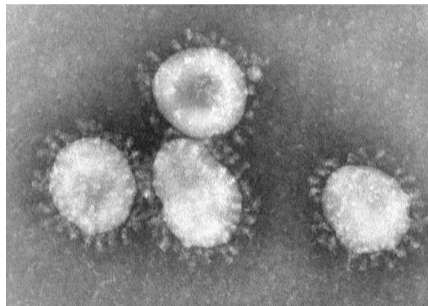- 293406 reads, 27 % virus, rest human
- Median read length 2.5 kb



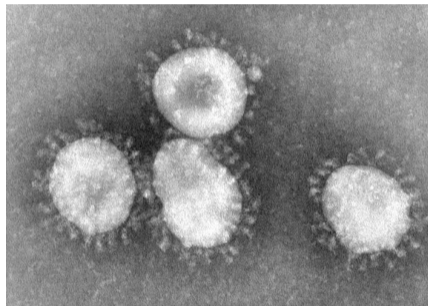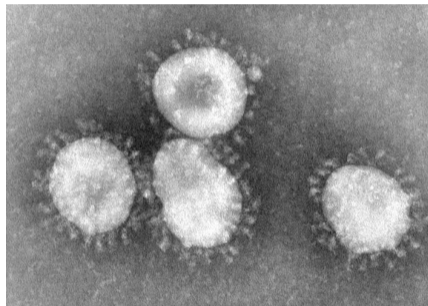en.wikipedia.org/wiki/Coronavirus#/media/

File:Coronaviruses_004_lores.jpg

# REAL CORONAVIRUS READ DATA

- HCoV 229E in human cell culture
- Direct RNA protocol kit on MinION
- 293406 reads, 27 % virus, rest human
- Median read length 2.5 kb
- Longest read: 26 kb (genome 27.3 kb)



`en.wikipedia.org/wiki/Coronavirus#/media/`

`File:Coronaviruses_004_lores.jpg`

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

# REAL CORONAVIRUS READ DATA

- HCoV 229E in human cell culture
- Direct RNA protocol kit on MinION
- 293406 reads, 27 % virus, rest human
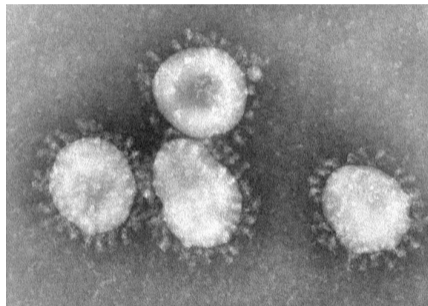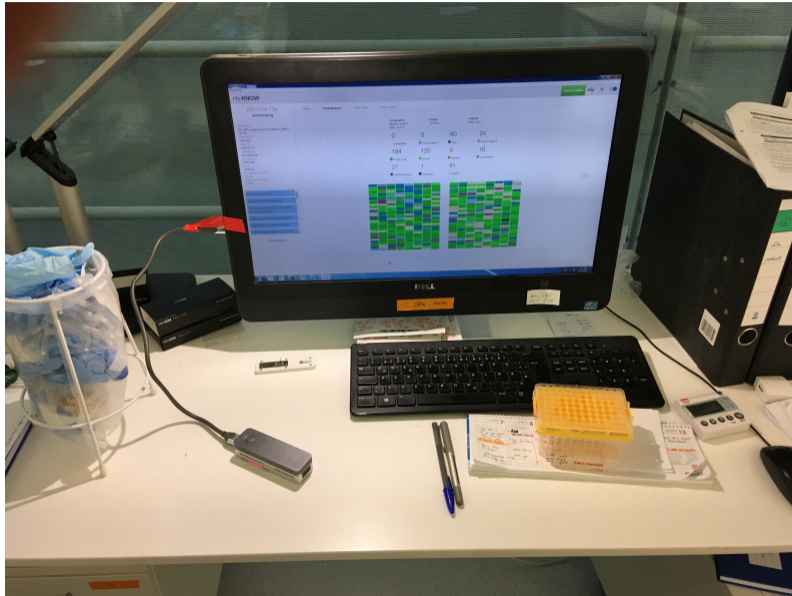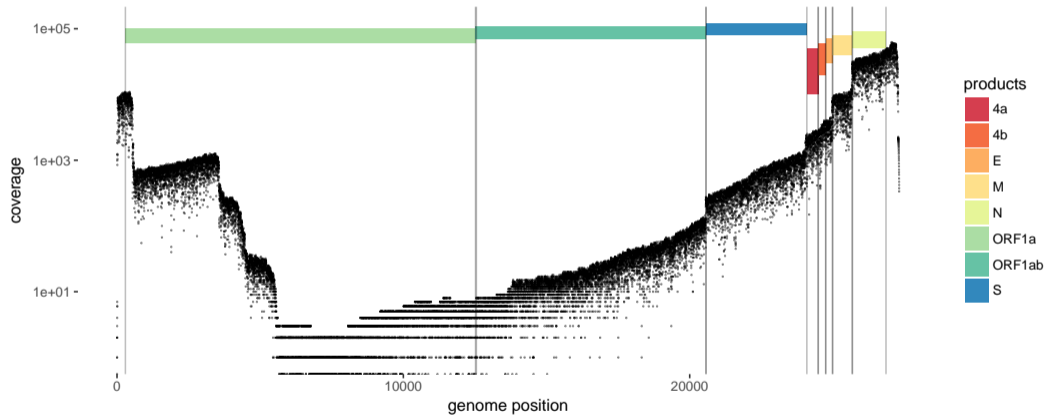- Median read length 2.5 kb
- Longest read: 26 kb (genome 27.3 kb)
- Error rate 15 % - mainly indels



en.wikipedia.org/wiki/Coronavirus#/media/

File:Coronaviruses_004_lores.jpg

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

Human Coronavirus 229E
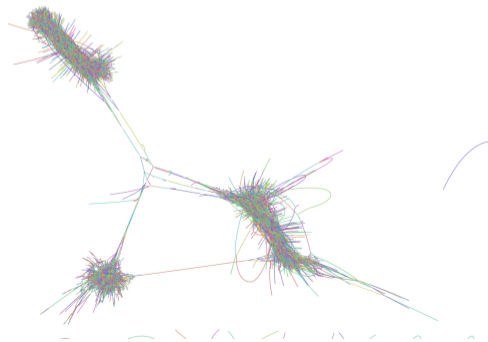
Subgenomic RNAs

corona - sequenced data - **2000 reads** - **k=40**

# SUBGRAPH CONSENSUS

- ▶ Needed: subgraph separation

# SUBGRAPH CONSENSUS

- Needed: subgraph separation

# SUBGRAPH CONSENSUS

- Needed: subgraph separation
- Implemented with min-cut

# SUBGRAPH CONSENSUS

- Needed: subgraph separation
- Implemented with min-cut
- Separates clusters that are minimally connected

## SUBGRAPH CONSENSUS

- Needed: subgraph separation
- Implemented with min-cut
- Separates clusters that are minimally connected
- Subgraph consensus is implemented

HCoV

S. cerevisiae

10086nt HCoV

HCoV

H sapiens

Yeast enolase is included in the direct RNA kit as a positive control

corona - sequenced data - **73533 reads** - **k=40**

# SEQUENCING ERRORS

- ▶ Mostly insertions and deletions

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

# SEQUENCING ERRORS

- ▶ Mostly insertions and deletions
- ▶ Of those: 70-80% deletions

# SEQUENCING ERRORS

- Mostly insertions and deletions
- Of those: 70-80% deletions
- Deletions happen systematically at homopolymers

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

# LONG READ ERROR CORRECTION

- Self-correction:
  Systematic errors are problematic

# LONG READ ERROR CORRECTION

- Self-correction:
  Systematic errors are problematic
- Hybrid correction with i.e. Illumina data:
  Alignment to noisy long reads difficult

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

# LONG READ ERROR CORRECTION

- ▶ Self-correction:
  Systematic errors are problematic
- ▶ Hybrid correction with i.e. Illumina data:
  Alignment to noisy long reads difficult
- ▶ `HG-CoLoR` by P. Morisse et al.

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

# LONG READ ERROR CORRECTION

- Self-correction:
  Systematic errors are problematic
- Hybrid correction with i.e. Illumina data:
  Alignment to noisy long reads difficult
- `HG-CoLoR` by P. Morisse et al.
- Longest read (25932 nt)
  Identity to reference: 84% → 99%
  Gap of 407 nt, 90 min runtime

FRIEDRICH-SCHILLER-
UNIVERSITÄT
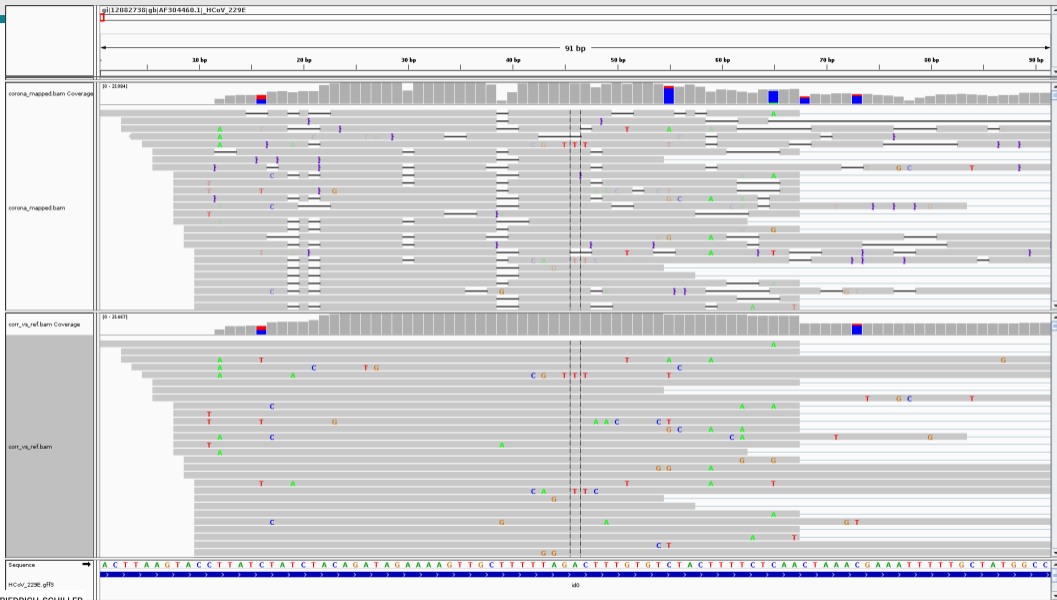JENA

# REFERENCE-BASED INDEL CORRECTION

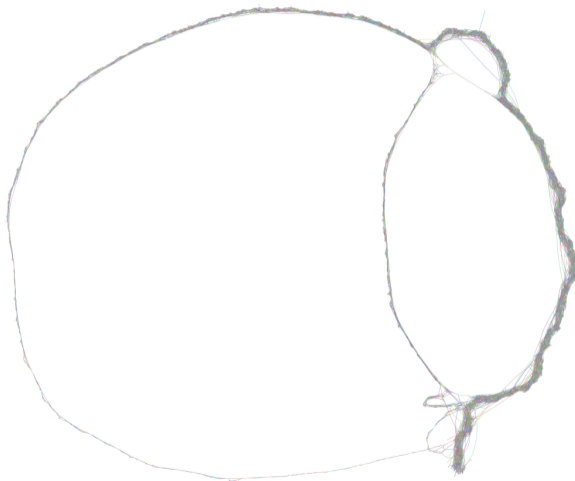- Reference from `nanopolish` by J. Simpson – RNA NYI

# REFERENCE-BASED INDEL CORRECTION

- ▶ Reference from `nanopolish` by J. Simpson – RNA NYI
- ▶ Align long reads to reference

# REFERENCE-BASED INDEL CORRECTION

- ▶ Reference from `nanopolish` by J. Simpson – RNA NYI
- ▶ Align long reads to reference
- ▶ Parse CIGAR string to remove insertions, fill deletions

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

Results so far
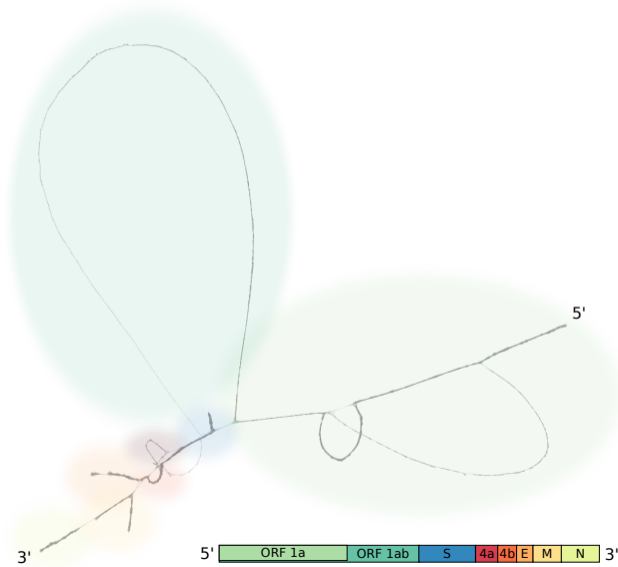
FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

corona - indels corrected - **1% best nucleotides** - **k=40**
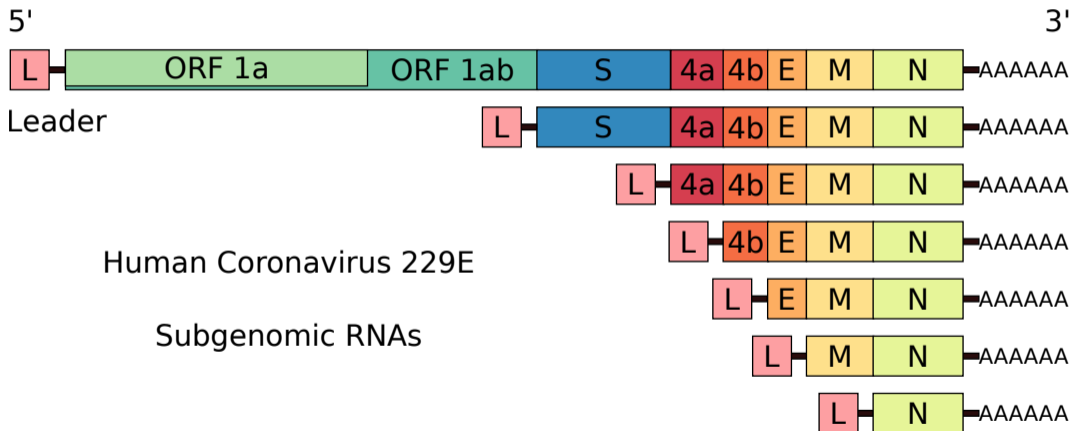
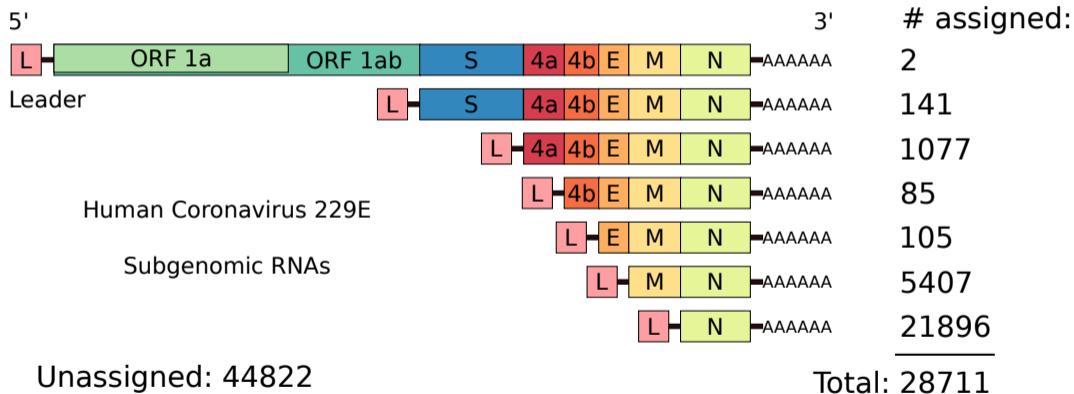corona - indels corrected - **10% best nucleotides** - **k=30**

corona - indels corrected - **1% best nucleotides** - **k=20**

5'

3'

5' | ORF 1a | ORF 1ab | S | 4a | 4b | E | M | N | 3'

## SUBGENOMIC TYPES



Human Coronavirus 229E
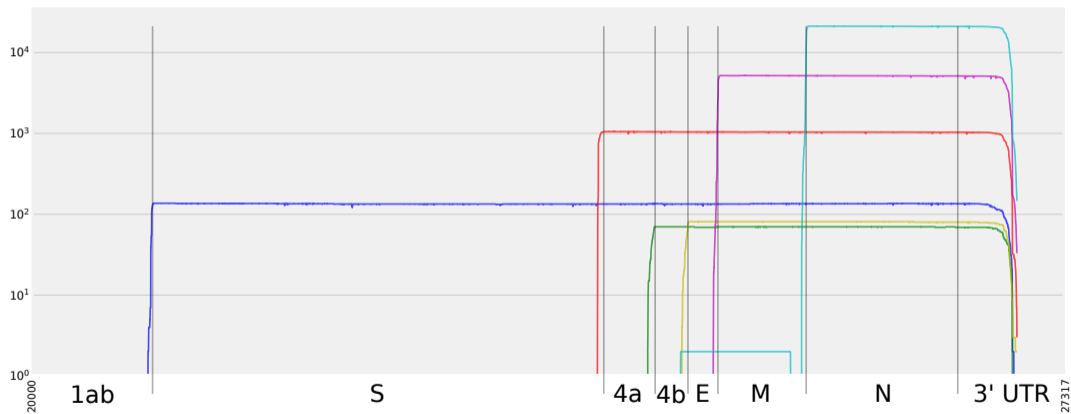
Subgenomic RNAs

# ANNOTATION BASED CLASSIFICATION

Coverage for subgenomic types
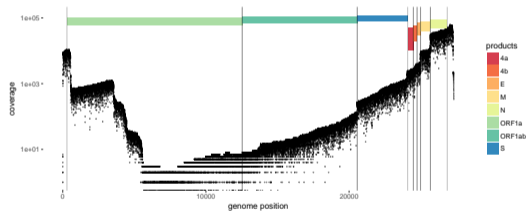
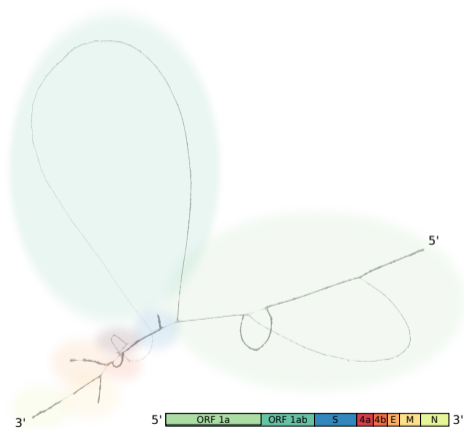■ S  ■ 4a  ■ 4b  ■ E  ■ M  ■ N

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

## CONCLUSIONS

- ▶ Viral full genome sequencing

## CONCLUSIONS

- ▶ Viral full genome sequencing
- ▶ Structure is visible in graph



5' | ORF 1a | ORF 1ab | S | 4a | 4b | E | M | N | 3'

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

# CONCLUSIONS
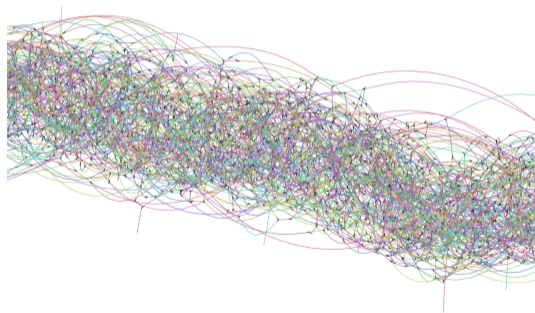
- ▶ Viral full genome sequencing
- ▶ Structure is visible in graph
- ▶ Importance of k

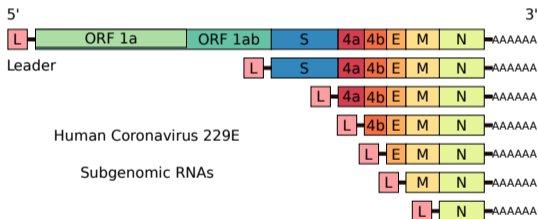## CONCLUSIONS

- ▶ Viral full genome sequencing
- ▶ Structure is visible in graph
- ▶ Importance of k
- ▶ Indel correction required

## CONCLUSIONS

- Viral full genome sequencing
- Structure is visible in graph
- Importance of k
- Indel correction required
- Coronavirus is ...complicated

# Outlook

## NEXT STEPS

- Test on new data (more HCoV, plum pox virus)

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

# NEXT STEPS

- Test on new data (more HCoV, plum pox virus)
- Improve error correction

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

## NEXT STEPS

- Test on new data (more HCoV, plum pox virus)
- Improve error correction
- Find robust way to extract haplotypes

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

## NEXT STEPS

- Test on new data (more HCoV, plum pox virus)
- Improve error correction
- Find robust way to extract haplotypes
- Utilize long read information

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

Andreas Goral

Adrian Viehweger
Celia Diezel

Manja Marz

Ramakanth Madhugiri
John Ziebuhr

All of my group!

**Thank you!**

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA