

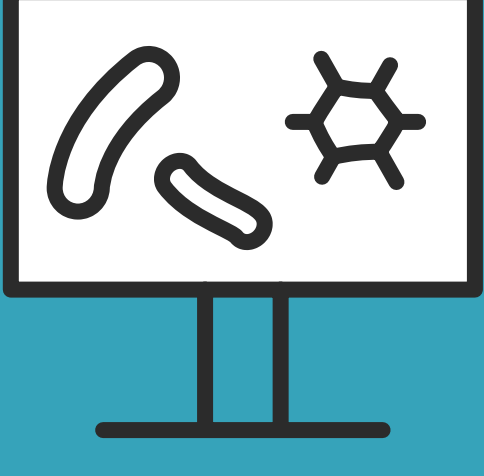
Benchmarking metagenomics classifiers on ancient viral DNA: a simulation study

BACKGROUND



Thanks to technological advances it is now possible to retrieve and sequence DNA from ancient samples. This genomic data includes DNA from ancient microbes.

AIM



Several classifiers have been developed to characterise the microbiota



Which classifier is best suited to screen ancient samples for ancient viruses?

Tested classifiers: Centrifuge¹, Kraken2², DIAMOND³ and MetaPhlAn2⁴

METHODS

1. Select reference genomes



All human DNA viruses stored in RefSeq (233)



Human genome reference

2. Simulate sequencing reads

simulated reads with ART⁵ and gargammel⁶



- Initial simulation: 60 bp reads
- Varying read length
- Adding deamination damage
- Adding sequencing errors

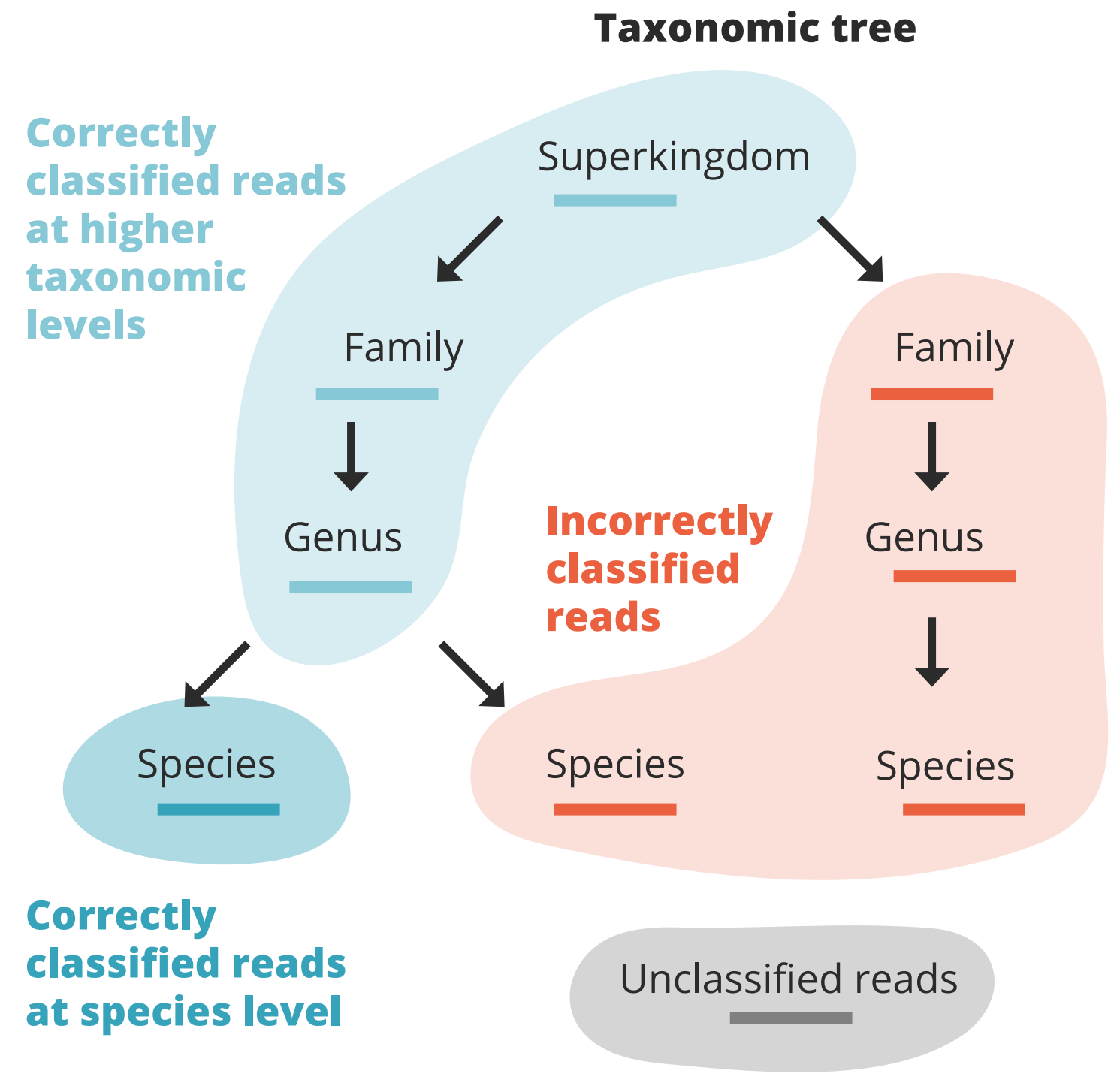
3. Read classification

Simulated reads



Centrifuge, Kraken2, DIAMOND, MetaPhlAn2

Classified reads

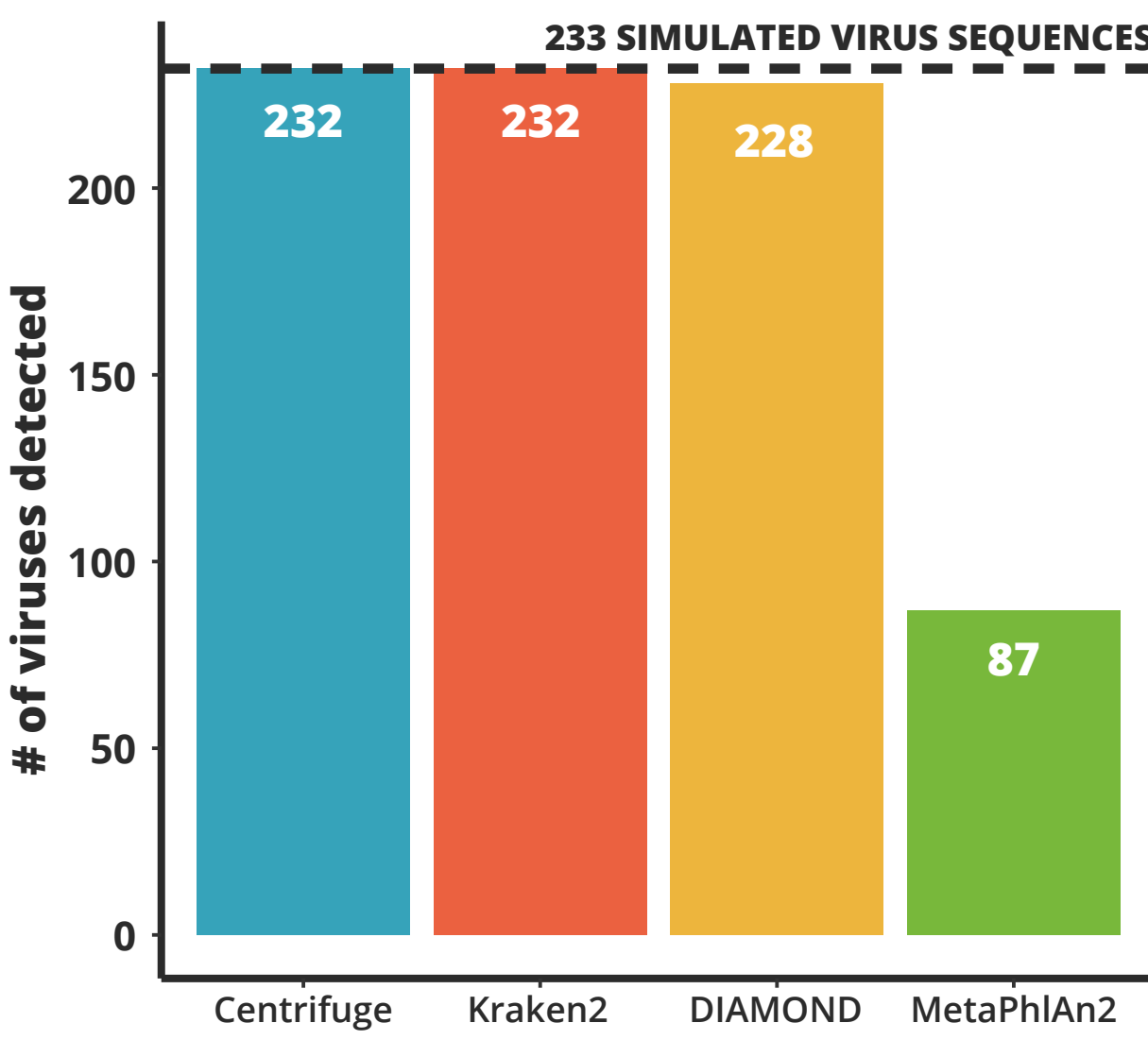


RESULTS

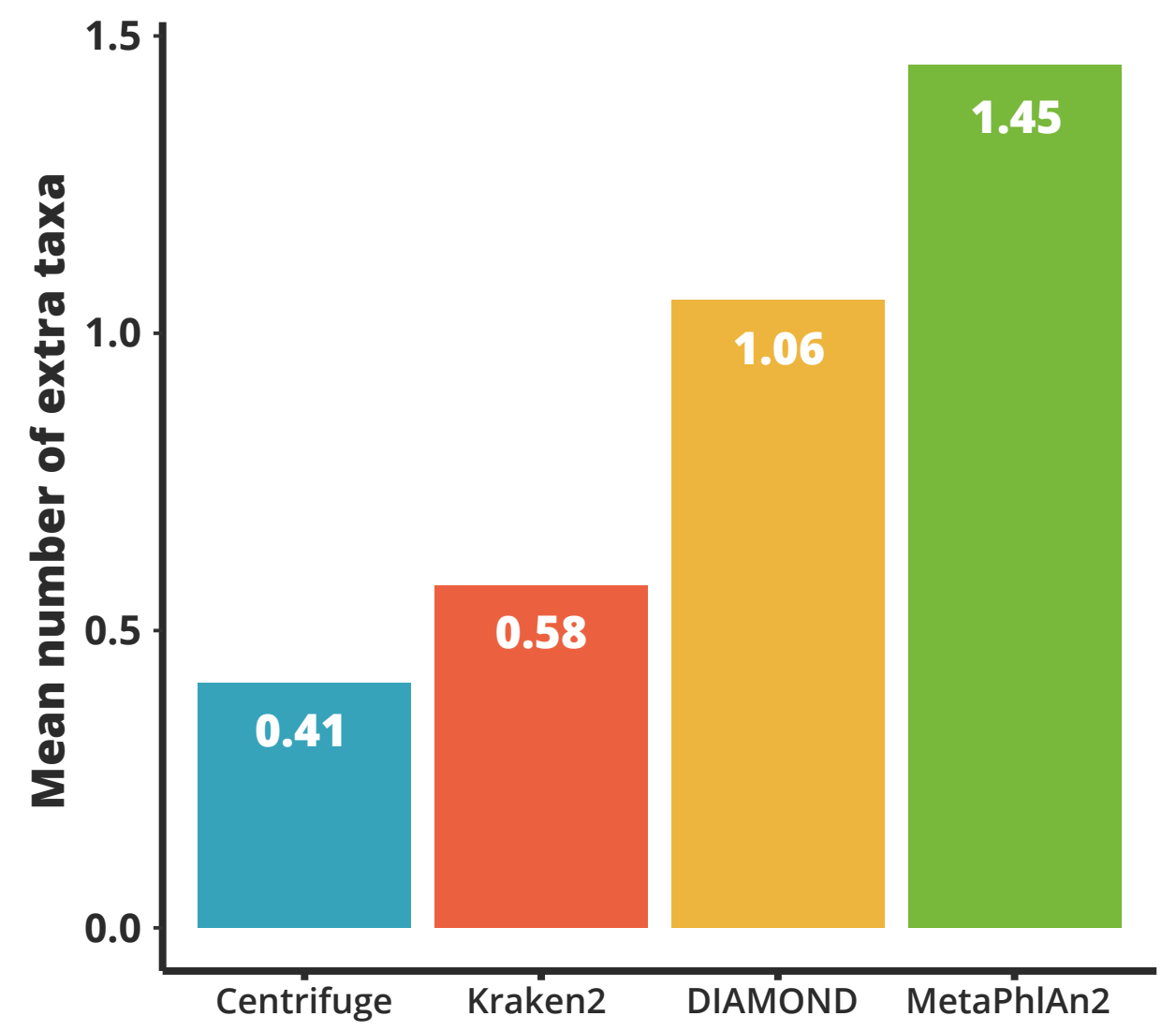
Compute precision and sensitivity

For 60 bp initial simulations sets:

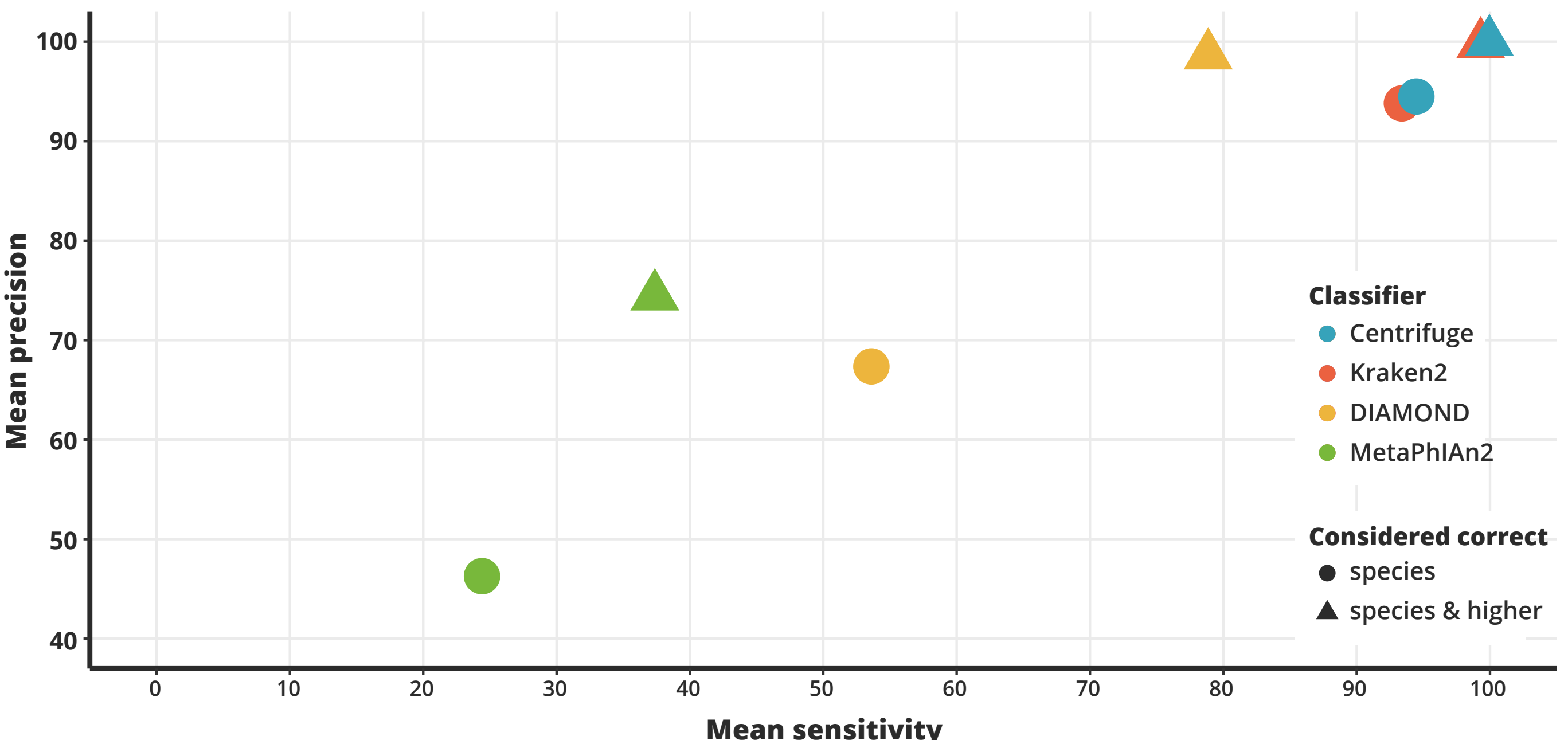
A Total number of correctly detected viruses



B Mean number of spurious extra taxa



C Sensitivity vs. precision



CONCLUSIONS

- Centrifuge** has:
 - The **highest sensitivity and precision** values
 - The best performance with **short (30 bp) reads**
- Most of the viruses were detected** by almost all the classifiers tested.
- Longer reads are better classified**, while increased sequencing error and increased deamination damage worsen the classification slightly.
- We recommend:**
 - using **strong filters** to remove human DNA
 - verifying that the **genomes of interest are included** in the classifiers' databases



- Centrifuge: Kim D, Song L, Breitwieser FP, Salzberg SL. 2016. Centrifuge: rapid and sensitive classification of metagenomic sequences. *Genome Research* 26:1721-1729 DOI 10.1101/gr.210641.116.
- Kraken2: Wood DE, Lu J, Langmead B. 2019. Improved metagenomic analysis with Kraken 2. *Genome Biology* 20:257 DOI 10.1186/s13059-019-1891-0.
- DIAMOND: Buchfink B, Xie C, Huson DH. 2015. Fast and sensitive protein alignment using DIAMOND. *Nature Methods* 12:59-60 DOI 10.1038/nmeth.3176.
- MetaPhlAn2: Truong DT, Franzosa EA, Tickle TL, Scholz M, Weingart G, Pasolunghi E, Tett A, Huttenhower C, Segata N. 2015. MetaPhlAn2 for enhanced metagenomic taxonomic profiling. *Nature Methods* 12:902-903 DOI 10.1038/nmeth.3589.
- ART: Huang W, Li L, Myers JR, Marth GT. 2012. ART: a next-generation sequencing read simulator. *Bioinformatics* 28:593-594 DOI 10.1093/bioinformatics/btr708.
- Gargammel: Renaud G, Hangh j K, Willerslev E, Orlando L. 2017. gargammel: a sequence simulator for ancient DNA. *Bioinformatics* 33:577-579.