# Part 3: Kraken2 Statistics

Kraken2 classification results have been compiled and saved as CSV files.

The Kraken2 results data includes:
– Sample names
– Species classifications
– Unaligned clade reads
– Non–human clade reads
– Dataset information (unaligned/nonhuman)
– Merged input read counts from RRstats
– Additional optional columns (confidence, minimum hit groups, etc.)

Data source: Combined from unaligned_kreports/ and nonhuman_kreports/
Processing: output_processing.R with species filtering and merging
Merging: Combined with runtime and read statistics by sample name
Output location: Check outputs/ directory for:
  – sample_report_data_with_metadata*op.csv (combined data)
  – sample_report_data_[dataset]_*op.csv (dataset–specific files)

Key findings:
– Average number of species found per sample before secondary filtering: 611.55
– Average number of species found per sample after secondary filtering: 80.075
– Number of unique species across all unaligned samples: 579

Homo sapiens: 66784243.6
Severe acute respiratory syndrome–related coronavirus: 3677261.7
Fusobacterium pseudoperiodonticum: 541689.5
Solobacterium moorei: 460470.5
Cupriavidus pauculus: 378178.25
Veillonella atypica: 293441.4
Sphingomonas paucimobilis: 263914.894736842
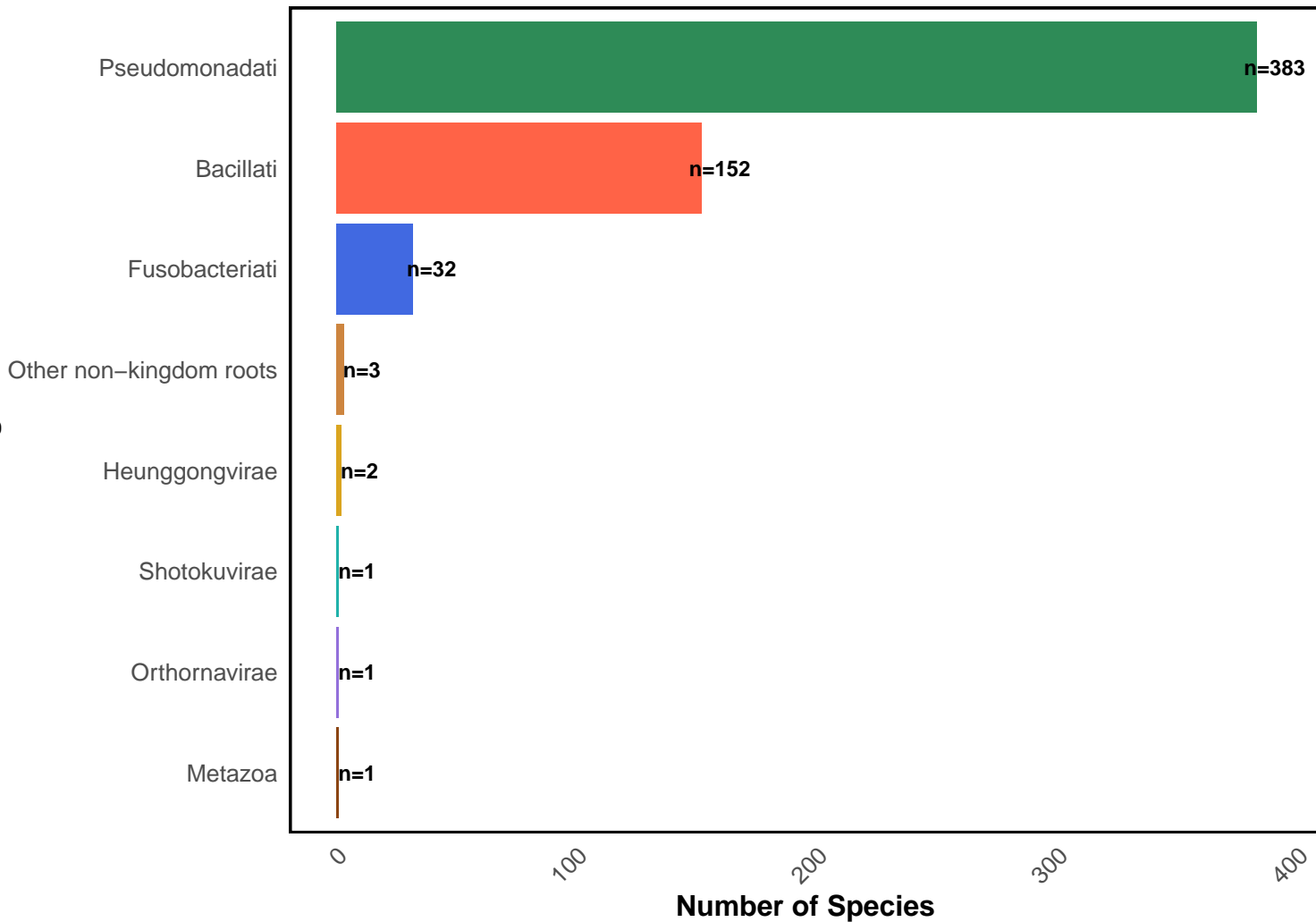Staphylococcus aureus: 223765.111111111
Selenomonas sputigena: 167271
Haemophilus parahaemolyticus: 162031

**Species per Kingdom – All Samples**

Total species: 575

Kingdom / Number of Species:

- Pseudomonadati — n=383
- Bacillati — n=152
- Fusobacteriati — n=32
- Other non–kingdom roots — n=3
- Heunggongvirae — n=2
- Shotokuvirae — n=1
- Orthornavirae — n=1
- Metazoa — n=1
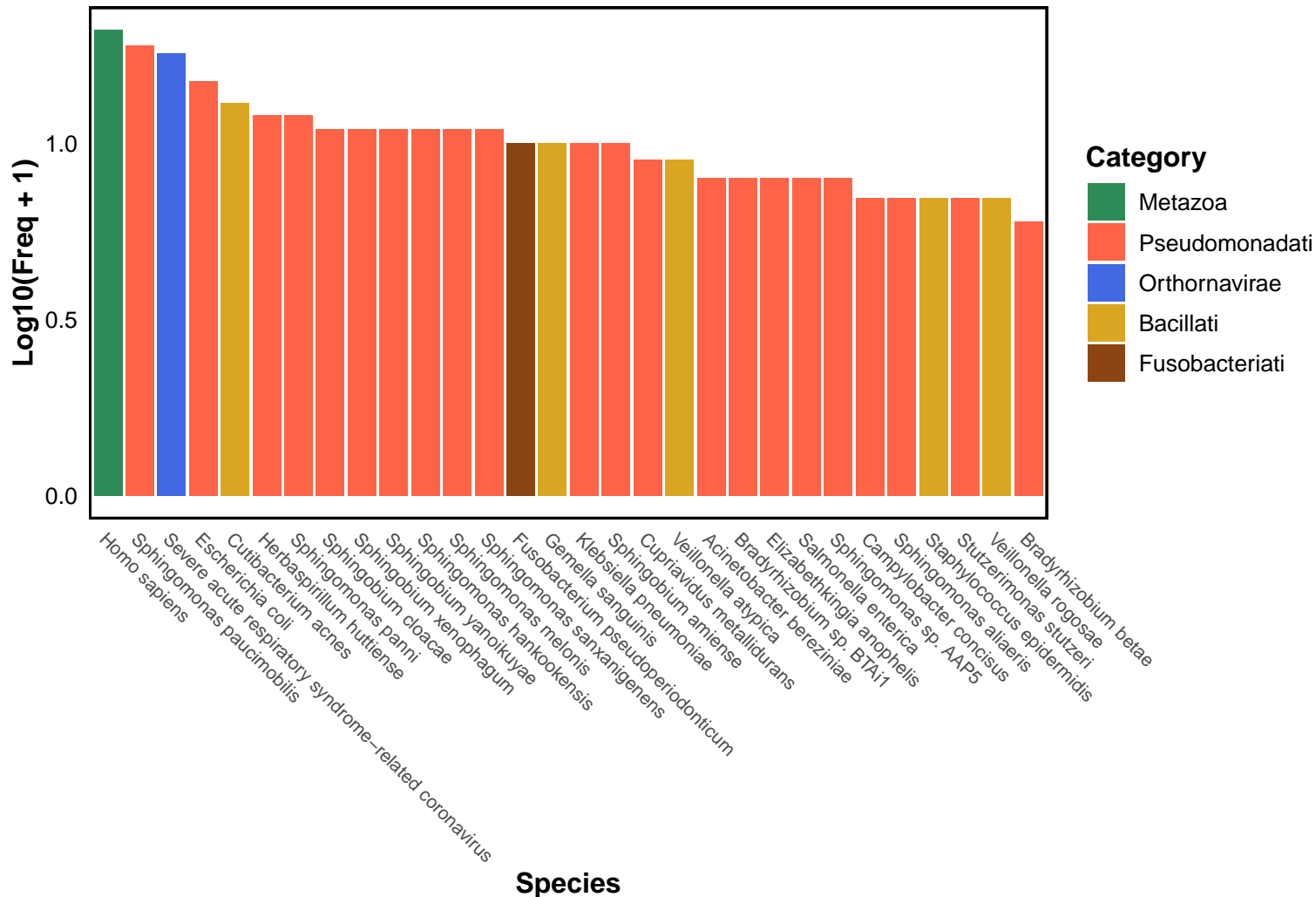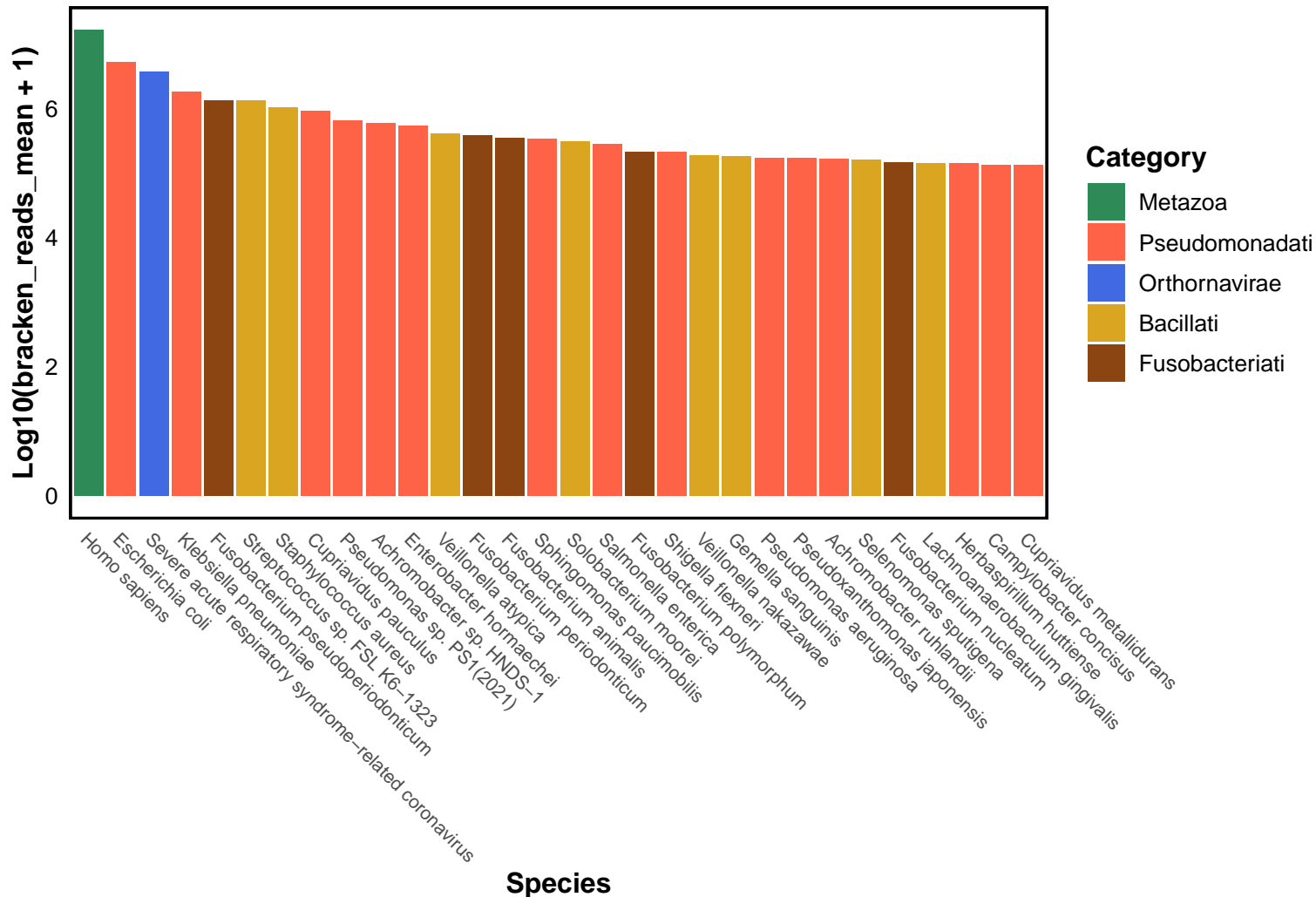
**Species by Kingdom – All Samples**

Top 30 species by Log10(Freq + 1)

# Species by Kingdom – All Samples

Top 30 species by Log10(bracken_reads_mean + 1)

**Y-axis:** Log10(bracken_reads_mean + 1)

**X-axis:** Species

**Category**
- Metazoa
- Pseudomonadati
- Orthornavirae
- Bacillati
- Fusobacteriati

Species (left to right):
Homo sapiens, Escherichia coli, Severe acute respiratory syndrome–related coronavirus, Klebsiella pneumoniae, Fusobacterium pseudoperiodonticum, Streptococcus sp. FSL K6–1323, Staphylococcus aureus, Cupriavidus pauculus, Pseudomonas aeruginosa, Achromobacter sp. PS1(2021), Enterobacter hormaechei, Veillonella atypica, Fusobacterium periodonticum, Fusobacterium animalis, Sphingomonas paucimobilis, Solobacterium moorei, Salmonella enterica, Fusobacterium flexneri, Shigella polymorphum, Veillonella nakazawae, Gemella sanguinis, Pseudomonas aeruginosa, Pseudoxanthomonas japonensis, Achromobacter ruhlandii, Selenomonas sputigena, Fusobacterium nucleatum, Lachnoanaerobaculum gingivalis, Herbaspirillum huttiense, Campylobacter concisus, Cupriavidus metallidurans

Kraken2 Clade Reads vs Bracken Reads (Top 25 Species)

Ranked by Kraken2 clade reads (log10 scale)

**Species by Kingdom – All Samples**

Top 30 species by avg_percentage

Category
- Pseudomonadati
- Orthornavirae
- Bacillati
- Fusobacteriati