Reconciliation and GMPD2 harmonization pipeline #1 (Gibb et al., 2021, biorxiv) HP3 Host and virus names reconciled to NCBI taxonomy **CLOVER** Temporal and sampling metadata are standardized Shaw Manual curation of unmatched names for consistency · Limited to mammals (v0.1.1) or EID2 vertebrates (v1.0.0) 2015 static **GLOBI** Metagenomic pipeline Every sample in the Sequence Read Archive NCBI is analyzed using the k-mer based method described in Katz et al. (2021) GenBank Out of the pool of samples with taxonomically-valid host metadata, the maximum number of k-mer hits for each Additional curation of known virus is recorded, to create a host-**PREDICT** viral higher taxonomy virus edgelist with an ad hoc scoring system for evidence (in progress) A threshold score is selected (maximizing the kappa statistic) based on presence or absence of associations in CLOVER as a **NCBI** Δ "ground truth"; associations are filtered to SRA those above the threshold score

Reconciliation and harmonization pipeline #2

- Dynamic datasets (△) are programmatically updated
- Host and virus names reconciled to NCBI taxonomy
- Temporal and sampling metadata are standardized

VIRION

- Manual curation of unmatched names to other sources (e.g., IUCN)
- Additional quality checks (e.g., phage and non-pathogenic viruses removed by taxonomy)

Edgelist.csv

Host, Virus, ID (join metadata by these fields)

TaxonomyHost.csv

Host, HostTaxID, HostNCBIResolved, HostGenus, HostFamily, HostOrder, HostClass, HostSynonyms

TaxonomyVirus.csv

Virus, VirusTaxID, VirusNCBIResolved, VirusGenus, VirusFamily, VirusOrder, VirusClass

Provenance.csv.gz

ID, Database, DatabaseVersion, ReferenceText, PMID

Detection.csv.gz

ID, DetectionMethod, DetectionOriginal, HostFlagID, VirusFlagContaminant, NCBIAccession

Temporal.csv.gz

ID, PublicationYear, ReleaseYear, ReleaseMonth, ReleaseDay, CollectionYear, CollectionMonth, CollectionDay