# A New Assembly-based Shotgun Data Analysis Pipeline for Microbiome Exploration in Online Analysis Platform, Nephele

**Angelina G Angelova, Duc Doan, Poorani Subramanian, Mariam Quiñones, Lewis Kim, Michael Dolan and Darrell E. Hurt**

All authors are affiliated with Bioinformatics and Computational Biosciences Branch, Office of Cyber Infrastructure and Computational Biology, National Institute of Allergy and Infectious Diseases, National Institutes of Health, Bethesda, MD 20892, USA

## ABSTRACT

Many researchers lack the computational training or resources to process large metagenomic datasets to the level optimal for extracting biological data. To address this problem, our Nephele2 team at NIAID, developed a new user-friendly command line-free pipeline, Whole metaGenome Sequence Assembly pipeline, version 2 (WGSA2) designed to bridge the computational gap between short shotgun reads and metagenomic assemblies. Nephele's WGSA2 pipeline can process metagenomic datasets from complex communities of various environments including human and animal-associated or environmental microbiomes, inclusive of prokaryotic, micro-eukaryotic and viral organisms. Overall, the WGSA2 pipeline allows the user to easily gain understanding of their metagenomic samples, with investment in time, effort or command line computation.

## INTRODUCTION

A plethora of amazing tools are available online (e.g. HUMAnN3.0, MG-R IDseq) that can perform taxonomic and/or functional classification on shotgun metagenomic reads. However, due to their short nature, shotgun often fail to provide relevant biological information. Long-read b metagenomic tools are also abundant online (e.g. GhostKoala, iPath, Meta NCBI BLAST) and can capture copious amounts of biological information require computational power and effort to obtain. To bridge the gap bet these two analytical strategies, the National Institute of Allergy and Infec Diseases (NIAID) has engaged their cloud-based microbiome analysis plat Nephele2, to provide new user-friendly command line-free Whole Metager Sequence Assembly-based pipeline, WGSA2. The pipeline is designed researcher's shortcut through the computational requirements of shotgun processing, read assembly and assembly evaluation, allowing for foc researcher's time and efforts towards project-specific downstream investiga (gene mining, community & statistical explorations, etc.).

## PIPELINE FEATURES

- Trim, filter, error correct raw shotgun reads (fastp)
- Decontamination against a choice database of organisms (Kraken2)
- Per-sample reference-free *de novo* assembly (metaSPAdes)
- Gene prediction & annotation (Prodigal, eggNOG-mapper2)
- Pathway inference based on database of choice (MinPATH)
- Taxonomic and metabolic community matrix per dataset
- Metagenome-assembled genomes (MetaBAT2, CheckM)
- Antimicrobial resistance peptide prediction (AMRFinderPlus)
- Sequences and abundance scores of assembled features (R)
- Exploratory community statistics and visualizations (R)

### Submit your Paired End WGSA job to Nephele:

**Job Details**

**Description of the job:** my_metagenomic_dataset

**Host Decontamination DB:**
- ✓ Human or Mouse DB
- Marine DB
- Mosquito DB
- Nematode DB

**Run additional trim & filter**

**Average Read Quality:** 10

**Minimum Read Length:** 60

**Trimming of 5':** 20

**Trimming of 3':** 15

**Output TEDread fastq files:** ☐

**Run AMRFinder:** ☑

**Metabolic pathways Database:** KEGG database
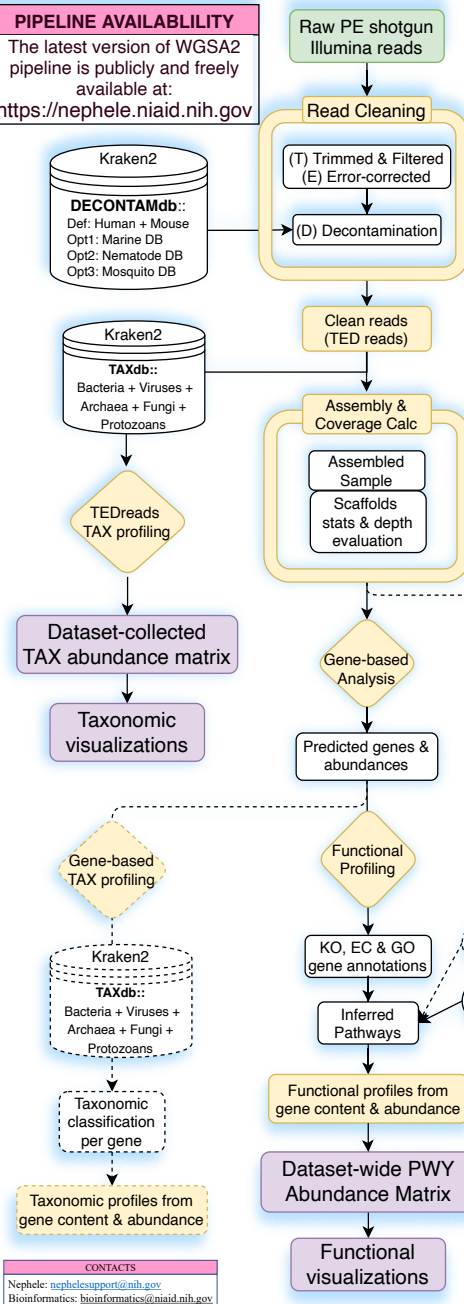
**Run Gene-based Taxonomic profiling:** ☑

**Produce taxonomic annotation on scaffolds:** ☐

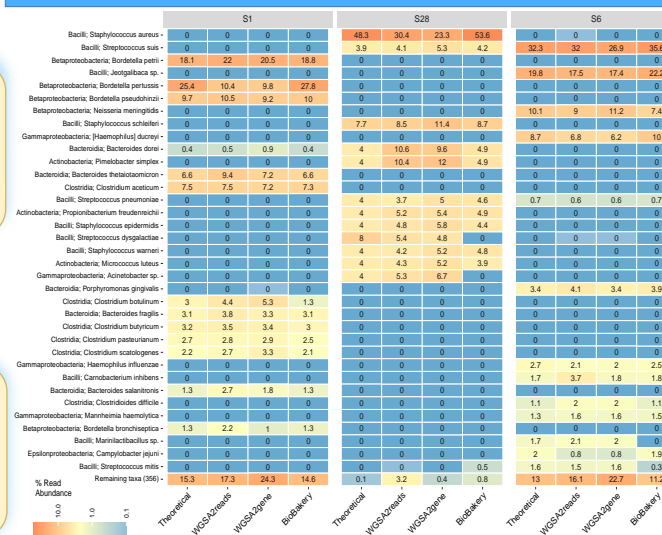**Produce MAGs (prokaryotic communities only):** ☑

## PIPELINE AVAILABILITY

The latest version of WGSA2 pipeline is publicly and freely available at:
https://nephele.niaid.nih.gov

### Pipeline flow

Raw PE shotgun Illumina reads → Read Cleaning

Kraken2
**DECONTAMdb::**
Def: Human + Mouse
Opt1: Marine DB
Opt2: Nematode DB
Opt3: Mosquito DB

Read Cleaning:
- (T) Trimmed & Filtered
- (E) Error-corrected
- (D) Decontamination

Clean reads (TED reads)

Kraken2
**TAXdb::**
Bacteria + Viruses + Archaea + Fungi + Protozoans

TEDreads TAX profiling

Assembly & Coverage Calc:
- Assembled Sample
- Scaffolds stats & depth evaluation

Dataset-collected TAX abundance matrix

Taxonomic visualizations

Gene-based TAX profiling

Kraken2
**TAXdb::**
Bacteria + Viruses + Archaea + Fungi + Protozoans

Taxonomic classification per gene

Taxonomic profiles from gene content & abundance

Gene-based Analysis → Predicted genes & abundances

Functional Profiling:
- KO, EC & GO gene annotations
- EC annots + MetaCyc maps
- KO annots + KEGG maps
- Inferred Pathways

MAGs-based Analysis → MAGs (Draft Genomes) → QC & TAX profiling
- Bin Taxonomy
- Bin Abundance
- Bin size (bp), Completeness, Contamination levels
- Bin Predicted Genes

Taxonomic profiles from draft genomes

Functional profiles from gene content & abundance

Dataset-wide PWY Abundance Matrix

Functional visualizations

### CONTACTS

Nephele: nephelesupport@nih.gov
Bioinformatics: bioinformatics@niaid.nih.gov

## BENCHMARKING



To assess the accuracy of the WGSA2 taxonomic profiling, the pipeline was run with 3 samples containing mock microbial communities from the 2nd CAMI Toy Human Microbiome Project Dataset (Sczyrba et al. 2017). The profiles were obtained using the Kraken2 classification tool against its standard database, and 2 of the available WGSA2 TAX profiling strategies - from cleaned short reads and from predicted genes (Prodigal), in assembled samples (metaSPAdes). These profiles were compared against each other, their theoretical counterparts and against another pipeline's shotgun-based profiles - BioBakery (McIver et al, 2018). Results indicate of closely related community profiles, composition and structure.

Poster # EEB1107

bioinformatics.niaid.nih.gov/metagenomics