# Alignment to Biothreats

Krista Ternus, PhD, PMP

Signature Science, LLC

# SeqScreen Publications

## Genome Biology

**SOFTWARE**

**Open Access**

## SeqScreen: accurate and sensitive functional screening of pathogenic sequences via ensemble learning

Advait Balaji[1†], Bryce Kille[1†], Anthony D. Kappell[2], Gene D. Godbold[3], Madeline Diep[4], R. A. Leo Elworth[1], Zhiqin Qian[1], Dreycey Albin[1], Daniel J. Nasko[5], Nidhi Shah[5], Mihai Pop[5], Santiago Segarra[6], Krista L. Ternus[2*] and Todd J. Treangen[1*]

[†]Advait Balaji and Bryce Kille contributed equally.

*Correspondence: kternus@signaturescience.com; treangen@rice.edu

[2]Signature Science, LLC, 8329 North Mopac Expressway, Austin, TX, USA
[1]Department of Computer Science, Rice University, Houston, TX, USA

**Abstract**

The COVID-19 pandemic has emphasized the importance of accurate detection of known and emerging pathogens. However, robust characterization of pathogenic sequences remains an open challenge. To address this need we developed SeqScreen, which accurately characterizes short nucleotide sequences using taxonomic and functional labels and a customized set of curated Functions of

## SeqScreen-Nano: a computational platform for streaming, in-field characterization of microbial pathogens

Advait Balaji
advait@rice.edu
Department of Computer Science,
Rice University
Houston, Texas, USA

Yunxi Liu
Department of Computer Science,
Rice University
Houston, Texas, USA

Michael G. Nute
Department of Computer Science,
Rice University
Houston, Texas, USA

Bingbing Hu
Department of Computer Science,
Rice University
Houston, Texas, USA

Anthony D. Kappell
Signature Science, LLC
Austin, Texas, USA

Danielle S. LeSassier
Signature Science, LLC
Austin, Texas, USA

Gene D. Godbold
Signature Science, LLC
Charlottesville, Virginia, USA

Krista L. Ternus
kternus@signaturescience.com
Signature Science, LLC
Austin, Texas, USA

Todd J. Treangen
treangen@rice.edu
Department of Computer Science,
Rice University
Houston, Texas, USA

**Abstract**

The COVID-19 pandemic forever underscored the need for bio-surveillance platforms capable of rapidly detecting emerging pathogens. Oxford Nanopore Technology (ONT) couples long-read sequencing with in-field capability, opening the door to real-time, in-field biosurveillance. Though a promising technology, streaming assignment of accurate functional and taxonomic labels with nanopore reads remains challenging given: (i) individual reads can span mul-

**CCS Concepts**

• Applied computing → Bioinformatics; Computational genomics.

**Keywords**

pathogen identification, metagenomics, bioinformatics

https://link.springer.com/article/10.1186/s13059-022-02695-x

https://dl.acm.org/doi/pdf/10.1145/3584371.3612960

# Publications about Sequences of Concern

## Categorizing Sequences of Concern by Function To Better Assess Mechanisms of Microbial Pathogenesis

Gene D. Godbold [a], Anthony D. Kappell [b], Danielle S. LeSassier [b], Todd J. Treangen [c], Krista L. Ternus [b]

[a] Signature Science, LLC, Charlottesville, Virginia, USA
[b] Signature Science, LLC, Austin, Texas, USA
[c] Department of Computer Science, Rice University, Houston Texas, USA

**ABSTRACT** To identify sequences with a role in microbial pathogenesis, we assessed the adequacy of their annotation by existing controlled vocabularies and sequence databases. Our goal was to regularize descriptions of microbial pathogenesis for improved integration with bioinformatic applications. Here, we review the challenges of annotating sequences for pathogenic activity. We relate the categorization of more than 2,750 sequences of pathogenic microbes through a controlled vocabulary called Functions of Sequences of Concern (FunSoCs). These allow for an ease of description by both humans and machines. We provide a subset of 220 fully annotated sequences in the supplemental material as examples. The use of this compact (~30 terms), controlled vocabulary has potential benefits for research in microbial genomics, public health, biosecurity, biosurveillance, and the characterization of new and emerging pathogens.

https://journals.asm.org/doi/epub/10.1128/iai.00334-21

Check for updates

## Improved understanding of biorisk for research involving microbial modification using annotated sequences of concern

Gene D. Godbold [1]*, F. Curtis Hewitt [2], Anthony D. Kappell [2], Matthew B. Scholz [2], Stacy L. Agar [1], Todd J. Treangen [3], Krista L. Ternus [2], Jonas B. Sandbrink [4] and Gregory D. Koblentz [5]*

[1] Signature Science, LLC, Charlottesville, VA, United States, [2] Signature Science, LLC, Austin, TX, United States, [3] Department of Computer Science, Rice University, Houston, TX, United States, [4] Nuffield Department of Medicine, University of Oxford, Oxford, United Kingdom, [5] Schar School of Policy and Government, George Mason University, Arlington, VA, United States

https://www.frontiersin.org/articles/10.3389/fbioe.2023.1124100/full

RICE    aclid    signature science LLC

# SeqScreen Installation

- SeqScreen documentation: https://gitlab.com/treangenlab/seqscreen/-/wikis/home
  - Runs in a Linux environment (needs up to 256 GB RAM, depending on use case)
  - We recommend installing SeqScreen via conda (or mamba)

- If you are not able to use conda, you can also download SeqScreen to run via Singularity or Docker thanks to quay.io
  - See SeqScreen versions here: https://quay.io/repository/biocontainers/seqscreen?tab=tags
  - `singularity pull docker://quay.io/biocontainers/seqscreen:4.2--hdfd78af_0`
  - `docker pull quay.io/biocontainers/seqscreen:4.2--hdfd78af_0`

- SeqScreen database, version released in March 2023:
  - https://s3.wasabisys.com/seqscreenv4/SeqScreenDB_23.3.tar.gz
  - md5sum = 4f01938a1f8d1a61e52ef9165e737824
  - Compressed database file is ~170 GB, uncompressed database and dependences are ~234 GB
  - Uncompress the database directory after download and leave subdirectory structure as is

# SeqScreen Outputs



Computer-friendly text file format that enables easy downstream automated parsing

Human-friendly HTML report that allows end users to interactively explore the results

RICE   ∧aclid   signature science LLC

# SeqScreen Code and Documentation on GitLab



https://gitlab.com/treangenlab/seqscreen

https://gitlab.com/treangenlab/seqscreen/-/wikis/Home

# Single Gene Sequence Outputs

- Input file = `single_gene_sequences/single_gene_seqs.fasta`

- Output files = `single_gene_sequences/single_gene_seqs_precomputed_results/`

  - `final_tsv_report_single_gene_seqs.xlsx`
  - `report_generation.zip`
  - `seqscreen_command.txt`
  - `single_gene_sequences_key.xlsx`

`single_gene_seqs_sensitive/report_generation/`
`single_gene_seqs_seqscreen_report_pathgo.tsv`

The *seqscreen_report_pathgo.tsv output file content was copied and pasted into an Excel file for ease of interpretation

RICE  aclid  signature science LLC