

aligner called SNAP that is about 1,000–10,000 times faster than BLAST. We combined SNAP with in-house developed software and other available open-source algorithms to create SURPI. SURPI can analyze MiSeq system runs in 10 minutes to a few hours, depending on the type of sequencing and analysis being performed. SURPI can also be run on a single computational server or even in the cloud. That's tractable in a clinical laboratory setting and was the motivation in creating SURPI.

“There are many patients where the pathogen can't be identified despite extensive conventional clinical testing.”

Q: How does SURPI identify pathogen DNA in the midst of all the metagenomic sequencing data?

GL: SNAP was written originally to align human genomes for resequencing applications in cancer and genetic testing. We modified and customized it for the extensive metagenomic databases provided by the National Center for Biotechnology Information (NCBI) that include sequences from pathogens and hosts. The SURPI pipeline uses SNAP to first identify and then subtract out human host sequences. It then aligns reads to reference sequences in NCBI databases, including GenBank, for comprehensive identification of all microorganisms: bacteria, viruses, fungi, and parasites. We added several programs to enable this data analysis to be performed in parallel.

Q: How long did it take you to develop SURPI?

GL: Developing the pipeline, and getting it validated, published, and made publicly available took about a year and a half.

Q: Why did you choose to use the MiSeq system to generate the metagenomic sequence for SURPI analysis?

GL: UCSF began using Illumina sequencing systems starting with the Genome Analyzer™ system and later upgraded to a HiSeq® system. When the MiSeq system became available, we migrated over to it for two reasons. It's a portable, desktop sequencer that's small enough to fit into a laboratory. The MiSeq system also generates a sizeable number of sequences (10–30 million) in 6–24 hours, a turnaround time that is suitable for infectious disease investigations.

Q: What library preparation kit did you use for this research study?

GL: We used TruSeq® kits, but recently switched to Illumina Nextera® library prep kits for this type of analysis. That's because Nextera kits are much faster.

Q: How did the MiSeq system perform in the study?

GL: I think the MiSeq system performed very well, producing 150 base pair (bp) reads that enabled us to make an unambiguous identification. Changes were made from the normal MiSeq system protocol so we could perform single-read sequencing, immediately pull the data from the instrument, and begin the analysis while the second paired end was being generated. By interrupting the sequencing run, we obtained data faster and that was critical. We're now working with Illumina to see if we can simultaneously analyze data in real time as it is generated by the instrument.

Q: What were the results of metagenomic sequencing and SURPI data analysis?

GL: In a cerebrospinal fluid (CSF) sample, we identified 475 of 3,063,784 sequence reads corresponding to *Leptospira* infection. This was a very convincing result because the sequences spanned the entire genome and *Leptospira* was the only credible pathogen DNA that was detected. When we sent the sample off for confirmatory testing to the CDC, it initially came back negative by antibody testing and the gold standard PCR for *Leptospira*. We later showed that the reason the CDC PCR results were negative was because the PCR assay they used had not been fully validated with that particular *Leptospira* species, *Leptospira santarosai*, and thus was not sensitive enough to detect the bacterial pathogen.

So it turned out that NGS was probably the only way the pathogen could have been identified, at least in an actionable timeframe. No clinical test available at the time would have been able to make the diagnosis, even if leptospirosis infection had been considered *a priori*.

Q: How much faster was the MiSeq system/SURPI method than current technologies?

GL: For this particular case, the sample-to-answer turnaround time was 48 hours. After the case report was published, we've decreased the turnaround time to less than 24 hours. We are ultimately aiming for an 8-hour turnaround time to make the NGS test fit within a single laboratory shift and be competitive with other molecular tests that are now performed routinely (multiplex PCR tests for virus detection, mass spectrometry-based methods for multiplex detection, and specific PCR-based tests to identify unusual organisms).

Q: What tests had been performed on this patient before enrolling him in your research study?

GL: The patient had been hospitalized three times over 4 months. Physicians had ordered an extensive infectious disease workup, including diagnostic laboratory testing, MRIs, CSF analysis, and a brain biopsy. All testing results were negative, although the spinal fluid and brain biopsy profiles showed prominent inflammation and strongly suggested the possibility of an infection.

