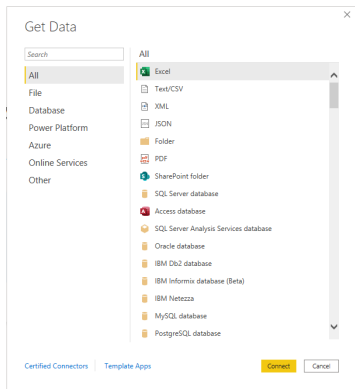


Colby Ford

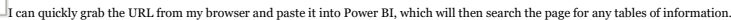
However, for any of us that are researchers in the bioinformatics and genomics space, we know that our files can be a bit difficult to work with. From FASTA to BAM, working with files in bioinformatics add a layer of uniqueness that requires some special care.

Today, if you take a look at Power BI Desktop's options for getting data, you'll see a ton of sources to which you can easily connect. One problem: none of these uniquely help us bioinformaticians.



In bioinformatics, there are a plethora of file types for every occasion. Among these are very popular ones such as FASTA (or FASTQ) and BAM and, more recently, GFF3 and BGEN. We can break these data sources down into three main types:

In Power BI, we can take advantage of Power Query to read in data and parse it appropriately. You'll notice that, while Power BI has tons of connects to everything from CSV files to Spark clusters, there are no built-in connectors for our beloved genomics file types (yet?). So, we'll have to use the Blank Query editor.



Navigator

Display Options ▾

HTML Tables [1]

Table 1

Suggested Tables [2]

Table 2

Table 3

Table ViewWeb View

Table 1

Column1	Column2	Column3	Column4
Name/Gene ID	Description	Location	Aliases
Select Item 43740578ORF1abID: 43740578	ORF1a polyprotein,ORF1ab polyprotein [Severe acute respiratory syndrom	NC_045512.2 (266..21555)	GU280_gp01
Select Item 43740577ORF8ID: 43740577	ORF8 protein [Severe acute respiratory syndrome coronavirus 2]	NC_045512.2 (27894..28259)	GU280_gp09
Select Item 43740576ORF10ID: 43740576	ORF10 protein [Severe acute respiratory syndrome coronavirus 2]	NC_045512.2 (29558..29674)	GU280_gp11
Select Item 43740575NID: 43740575	nucleocapsid phosphoprotein [Severe acute respiratory syndrome corona	NC_045512.2 (28274..29533)	GU280_gp10
Select Item 43740574ORF7bID: 43740574	ORF7b [Severe acute respiratory syndrome coronavirus 2]	NC_045512.2 (27756..27887)	GU280_gp08
Select Item 43740573ORF7aID: 43740573	ORF7a protein [Severe acute respiratory syndrome coronavirus 2]	NC_045512.2 (27394..27759)	GU280_gp07
Select Item 43740572ORF6ID: 43740572	ORF6 protein [Severe acute respiratory syndrome coronavirus 2]	NC_045512.2 (27202..27387)	GU280_gp06
Select Item 43740571MID: 43740571	membrane glycoprotein [Severe acute respiratory syndrome coronavirus;	NC_045512.2 (26523..27191)	GU280_gp05
Select Item 43740570EID: 43740570	envelope protein [Severe acute respiratory syndrome coronavirus 2]	NC_045512.2 (26245..26472)	GU280_gp04
Select Item 43740569ORF3aID: 43740569	ORF3a protein [Severe acute respiratory syndrome coronavirus 2]	NC_045512.2 (25393..26220)	GU280_gp03
Select Item 43740568SID: 43740568	surface glycoprotein [Severe acute respiratory syndrome coronavirus 2]	NC_045512.2 (21563..25384)	GU280_gp02, spike glycoprotein

Add Table Using Examples

LoadTransform DataCancel

This enables users to take advantage of data from virtually any site. Try it out on the [Protein Data Bank](#), [NCBI](#), [PlasmoDB](#), and more!

Takeaway Messages

- Be mindful of memory. Bioinformatics files can be large and, if you’re running on a machine with limited resources, you might bog it down.
- Check the defined specifications of any file format you’re looking to parse.
- R or Python can be your BFF, especially for binary or really complex file types.

Demo Video

Resources

All code used in above demos and additional examples are available at: www.github.com/BlueGranite/bioPowerBI

If you’d like to learn more about Power BI and how it can help you, [contact BlueGranite](#) today.