

SRA Handbook

Last Updated: January 14, 2016



National Center for Biotechnology Information (US)
Bethesda (MD)

National Center for Biotechnology Information (US), Bethesda (MD)

NLM Citation: SRA Handbook [Internet]. Bethesda (MD): National Center for Biotechnology Information (US); 2010-.

This documentation provides an overview and help manual for the Sequence Read Archive (SRA) at the National Center for Biotechnology Information.

Table of Contents

Using the SRA	1
Download Guide	3
Overview.....	3
Download with Prefetch.....	4
The Run Browser.....	5
SRA BLAST.....	7
Direct downloading of fasta and fastq format.....	9
Downloading metadata associated with SRA data files.....	9
Aspera Transfer Guide	13
Notice.....	13
Overview.....	13
Aspera.....	13
Using ascp to Download by Command Line.....	14
Using ascp to Upload by Command Line.....	15
Requirements.....	16
Troubleshooting.....	16

Using the SRA

Download Guide

Created: September 9, 2009; Updated: January 14, 2016.

Overview

The purpose of this document is to explain to users how to download datasets of interest and associated metadata.

Important Notes on Download Facilities

- One basic format (.sra) is provided by the SRA for all publicly available data. The SRA Toolkit is provided to allow conversion to several popular formats.
- At a minimum, users are advised to use Aspera Connect (or the equivalent command line tool, ascp) for bulk downloads, rather than HTTP or FTP. Aspera provides faster bandwidth, a higher level of flow control, user level encryption, and the ability to download trees of components.
- We most strongly recommend the use of the [SRA Toolkit](#) to download data files directly. The individual utilities are able to resolve SRA accessions and initiate downloads automatically. The ‘[prefetch](#)’ utility is specifically provided for researchers that wish to download SRA data using a command line utility.

Related Documents

[NCBI Large Data Download Best Practices](#)

Notices

Reference herein to any specific commercial products, process, or service by trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government, and shall not be used for advertising or product endorsement purposes.

Software Version

This guide is current to SRA Toolkit version 2.5.x. Instructions for previous versions of the SRA Toolkit may be different from those provided in this guide. We recommend that users stay current with [SRA Toolkit](#) updates to benefit from feature additions and bug fixes.

Reference Compression

Compression by reference is a sequence alignment compression process for storing data. Compression by reference stores the difference in base pairs between sequence data and the segment(s) to which it is aligned. Throughout this document you will note that the behavior and properties of reference compressed SRA data and conventional data differ significantly. Notably,

- The SRA Toolkit can output reference-compressed data as aligned sam and can perform pileup analysis.
- The SRA Toolkit requires internet connectivity in order to download reference sequences in order to process aligned SRA data.
- Only aligned data can be viewed in the NCBI Sequence Viewer.
- Aligned data cannot be filtered in the SRA Run Browser.
- Aligned data cannot currently be searched in SRA BLAST (this is actively being developed).

Download with Prefetch

The SRA Toolkit can be used to directly download SRA data files and reference sequences (see the “Reference Compression” section above). We strongly encourage users to use these methods to access SRA data as they are simple to use and they avoid many of the manual steps required by other methods (searching FTP directories, browsing and clicking, etc.).

The [SRA Toolkit](#) will have to be properly configured in order to access NCBI servers and download data. Recent versions of the Toolkit are packaged with a ‘default’ configuration that should work for most users. Please review the pros and cons for using the default configuration [here](#). If the default configuration does not work for your installation, or you wish to customize aspects of file handling by the Toolkit (e.g., where downloaded files are stored locally), you will need to [configure the Toolkit](#) and then [test it](#) to confirm that it is operating as expected. Please email sra@ncbi.nlm.nih.gov if you have any problems configuring or using the Toolkit.

Prefetch

The ‘[prefetch](#)’ utility in the SRA Toolkit can be used to download SRA data and any required reference sequences in a single operation. Prefetch can use either HTTP (default) or ascp (if installed) to contact the SRA, resolve the accessions that you have specified, and then download the data. Prefetch can be used on single data files or to batch download several at a time. Below is an example prefetch command with the expected output. More information can be obtained on the [prefetch documentation](#) page and by executing ‘prefetch --help’.

```
$ prefetch SRR390728
Maximum file size download limit is 20,971,520KB
2016-01-14T16:57:02 prefetch.2.5.7: 1) Downloading 'SRR390728'...
2016-01-14T16:57:02 prefetch.2.5.7: Downloading via fasp...
2016-01-14T16:57:08 prefetch.2.5.7: fasp download succeed
2016-01-14T16:57:08 prefetch.2.5.7: 1) 'SRR390728' was downloaded successfully
2016-01-14T16:57:09 prefetch.2.5.7: 'SRR390728' has 25 unresolved dependencies
2016-01-14T16:57:09 prefetch.2.5.7: 2) Downloading 'ncbi-acc:GPC_000000394.1?vdb-ctx=refseq'...
2016-01-14T16:57:09 prefetch.2.5.7: Downloading via fasp...
2016-01-14T16:57:13 prefetch.2.5.7: fasp download succeed
2016-01-14T16:57:13 prefetch.2.5.7: 2) 'ncbi-acc:GPC_000000394.1?vdb-ctx=refseq' was downloaded successfully
2016-01-14T16:57:13 prefetch.2.5.7: 3) Downloading 'ncbi-acc:GPC_000000395.1?vdb-ctx=refseq'...
2016-01-14T16:57:13 prefetch.2.5.7: Downloading via fasp...
2016-01-14T16:57:15 prefetch.2.5.7: fasp download succeed
```

Note that the example file is reference-compressed and that prefetch automatically obtains the reference sequences required to extract data from the .sra file. If your Toolkit installation is not properly configured, or you elect to block the ability of the Toolkit to contact NCBI, you will then need to determine (1) if your downloaded dataset is reference-compressed, (2) if so, which references are required to access the data (see [vdb-dump](#) for an example of how to determine this), and (3) acquire the reference sequences manually [here](#).

Other Toolkit utilities

All SRA Toolkit functions - most notably the ‘dump’ utilities that convert SRA data into other formats - are able to download data “on-the-fly” at runtime. This works like prefetch, as the tools will also automatically acquire all needed reference sequences. To invoke a Toolkit utility to download data as they are converted to your preferred format, simply execute the utility on an SRA accession rather than a local file. In other words, the command

```
$ fastq-dump --split-files SRR390728
```


Is implicitly requesting that fastq-dump download SRR390728 and its references from the SRA and then output the data in fastq format. Conversely,

```
$ fastq-dump --split-files ~/Downloads/SRA/SRR390728.sra
```

Is instructing fastq-dump to operate on a local file that was previously downloaded from the SRA. In this case fastq-dump would still attempt to contact NCBI to obtain the references needed to convert the data to fastq (unless you have specifically configured the Toolkit to not contact NCBI).

The Run Browser

The [SRA Run Browser](#) can display sequencing and instrumentation data on a given run. Typically the Run Browser is reached as a click through from an Entrez SRA Experiment report. Users may also navigate by entering a run accession (SRR, DRR, or ERR) directly in the Run Browser.

The screenshot shows the NCBI SRA Run Browser interface for run SRR390728. The top navigation bar includes links for Main, Browse, Search, Download, Submit, Documentation, Software, Trace Archive, Trace Assembly, Trace Home, and Trace BLAST. The main header indicates the run is an RNA-Seq (polyA+) analysis of DLBCL cell line HS0798. Below this, there are tabs for Metadata, Alignment, Reads, and Download. The Metadata tab is active, displaying a table with columns: Run, Spots, Bases, Size, GC content, Published, and Access Type. The table shows data for SRR390728. Below the table, there is a section for Reference, Length, Coverage, and Unaligned, showing data for NCBI36_BCCAGSC_variant. A quality graph is also visible. The bottom section shows the Experiment and Library details, including the Name, Platform, Strategy, Source, Selection, and Layout. The Biosample section provides a description of the sample and links to the Organism and Links. The Bioproject section shows the SRA Study and Title.

Run	Spots	Bases	Size	GC content	Published	Access Type
SRR390728	7.2 M	516.9 Mbp	193.6 M	45.7%	2011-12-21	public

Reference	Length	Coverage	Unaligned
NCBI36_BCCAGSC_variant	3.1 Gbp	0.15x	9.60%

Experiment	Library												
SRX079566	<table border="1"> <thead> <tr> <th>Name</th> <th>Platform</th> <th>Strategy</th> <th>Source</th> <th>Selection</th> <th>Layout</th> </tr> </thead> <tbody> <tr> <td>HS0798</td> <td>Illumina</td> <td>RNA-Seq</td> <td>TRANSCRIPTOMIC</td> <td>cDNA</td> <td>PAIRED</td> </tr> </tbody> </table>	Name	Platform	Strategy	Source	Selection	Layout	HS0798	Illumina	RNA-Seq	TRANSCRIPTOMIC	cDNA	PAIRED
Name	Platform	Strategy	Source	Selection	Layout								
HS0798	Illumina	RNA-Seq	TRANSCRIPTOMIC	cDNA	PAIRED								

Biosample	Sample Description	Organism	Links
SAMN00630374 (SRS212581)	established from ascites of a 45-year-old Caucasian man with diffuse large cell lymphoma	Homo sapiens	<ul style="list-style-type: none"> human B cell lymphoma cell line DB DSMZ:ACC-539

Bioproject	SRA Study	Title
PRJNA74797	SRP001599	The Cancer Genome Characterization Initiative: Parent Study

Viewing data in the Run Browser

Reference compressed (aligned) SRA data have an “Alignment” tab. Clicking on this tab will allow you to configure the [NCBI Sequence Viewer](#) to display the data aligned to a reference sequence.

NCBI Site map All databases Search

Sequence Read Archive

Main Browse Search Download Submit Documentation Software Trace Archive Trace Assembly Trace Home Trace BLAST

Studies Samples Analyses **Run Browser** Provisional SRA

RNA-Seq (polyA+) analysis of DLBCL cell line HS0798 (SRR390728) [Change accession...](#)

Metadata Alignment Reads Download

Alignment	Reads	Bases	Fraction
Primary	13.0 M	467.2 Mbp	90.4%

Reference

1 **Range** 1-1000000

[Homo sapiens chromosome 1, reference assembly, complete sequence](#)

[What does it do?](#)

View

scope	accession	count	in
<input checked="" type="radio"/> this run	SRR390728	1	Sequence Viewer
<input type="radio"/> same experiment	SRX079566	1	
<input type="radio"/> same sample	SRS212581	1	
<input type="radio"/> same study	SRP001599	10	
all sra		4,449	

Output this run in FASTA format to Screen File

To view the raw reads in a single Run, click on the “Reads” tab. Individual reads can be viewed and searched (see next section – note that only unaligned data can currently be searched). Various options can be applied using the “View” menu (e.g., display decimal quality scores, technical reads, etc.).

NCBI Site map All databases Search

Sequence Read Archive

Main Browse Search Download Submit Documentation Software Trace Archive Trace Assembly Trace Home Trace BLAST

Studies Samples Analyses **Run Browser** Provisional SRA

RNA-Seq (polyA+) analysis of DLBCL cell line HS0798 (SRR390728) [Change accession...](#)

Metadata Alignment **Reads** Download

< 1 1 717858 >

View: ☒ biological reads ☐ technical reads ☐ quality scores [advanced options](#)

Reads (separated)

1. [SRR390728.1](#) [SRS212581](#)
name: 1, member: default

2. [SRR390728.2](#) [SRS212581](#)
name: 2, member: default

3. [SRR390728.3](#) [SRS212581](#)
name: 3, member: default

4. [SRR390728.4](#) [SRS212581](#)
name: 4, member: default

5. [SRR390728.5](#) [SRS212581](#)
name: 5, member: default

6. [SRR390728.6](#) [SRS212581](#)
name: 6, member: default

7. [SRR390728.7](#) [SRS212581](#)
name: 7, member: default

8. [SRR390728.8](#) [SRS212581](#)
name: 8, member: default

9. [SRR390728.9](#) [SRS212581](#)
name: 9, member: default

10. [SRR390728.10](#) [SRS212581](#)
name: 10, member: default

```
>gnl|SRA|SRR390728.1.1 1 undefined (Biological, Reverse)
CATCTTCACGTAGTCTCGAGCCTGGTTTCAGC
>gnl|SRA|SRR390728.1.2 1 undefined (Biological, Reverse)
GATGGAGAAATGACTTTGACAACTGAGAGAGNTNC
```

The Run Browser supports IUPAC single letter nucleotide codes (data submitted in color space are presented in base space; the SRA Toolkit can be used to download and output the data in color space, if required). Quality scores are presented in the Phred scale.

Filtering and Selection

The Reads tab in the Run Browser can be used to filter and search reads according to certain regular expression pattern matching:

- Sequence substring: one of the biological reads for a spot should contain the substring.
Examples: `ATTGGA`, `^ATTGGA`, `ATTGGA$`, `ATGDNNAT`, and `ATGGA&GCGC`. The strings are case

insensitive, and belong to either 2NA or 4NA alphabets. String length limited to 29 characters in 4NA alphabet (includes IUPAC substitution codes) or 61 characters in 2NA alphabet (ACGT only). Search is case insensitive and strings may be combined with boolean operators & | ! (AND, OR, NOT). See "[SRA nucleotide search expressions](#)" for more details.

- Name of a spot you are looking for. Example: [EXWA4RL02G9Z6H](#)
- Name of sample pool member, or "all" for all members. Example: [M22_V2](#) will return all spots assigned to the sample pool member M22_V2 for run SRR031989.
- Spot Id. Example: [23](#)

Please note that the filter searches across read boundaries within each spot. Thus, pattern matches within technical reads and across paired-end data boundaries will also be returned.

The filter provided in the Run Browser has limited functionality, but is quite fast if you are looking to quickly search a single Run for a defined sequence of interest. Please see the section below on SRA BLAST if you require more advanced searching or searches across multiple sequencing libraries.

Downloading Data from the Run Browser

Clicking on the "Download" tab in the Run Browser will present a selection of links that will allow you to download (1) an individual dataset (Run), (2) all datasets in a given sequencing library (Experiment), or (3) all datasets linked to a given project (Study). You are also provided with three download choices: Aspera (using the [Aspera Connect plugin](#)), HTTP (using your browser), FTP (using command line FTP or a client).

NCBI Site map All databases Search

Sequence Read Archive

Main Browse Search Download Submit Documentation Software Trace Archive Trace Assembly Trace Home Trace BLAST

Studies Samples Analyses Run Browser Provisional SRA

RNA-Seq (polyA+) analysis of DLBCL cell line HS0798 (SRR390728) [Change accession...](#)

Metadata Alignment Reads Download

Object	.sra
Run SRR390728	193.6 Mb HTTP FTP Aspera
Experiment SRX079566	1.2 Gb HTTP FTP Aspera
Study SRP001599	14.0 Gb HTTP FTP Aspera

SRA BLAST

SRA BLAST can be used to for advanced searching of single or multiple sequencing libraries from the same or different projects. There are two ways to access SRA BLAST in order to build a "search space" from which you are attempting to pull matches to your sequence(s) of interest. Successful BLAST searches will lead you to a results / summary page that can be used to download reads of interest or be directed to the [SRA Run Browser](#) to further investigate or download the entire dataset.

Note that SRA BLAST currently has a limit of 2^{11} reads (approximately 2 billion) per search – attempts to add more than 2^{11} reads will result in an error and rejection of the search. Users that require more substantial search capability are advised to contact the SRA (sra@ncbi.nlm.nih.gov) to determine if other SRA BLAST tools might be of use.

Sending Entrez results to SRA BLAST

After performing an Entrez query to restrict results to datasets of interest, you may use the "Send to" feature to select datasets of interest and send them to SRA BLAST.

SRA search interface showing a query: `human[organism] NOT sra_nucore_alignment[Filter]`. The search results are displayed in a table with columns for Accession, Study Name, and Description. A 'Send to' dropdown menu is open, showing options: File, Clipboard, and BLAST. The results list includes studies like 'Cystic Fibrosis Longitudinal Study' and 'GSM1335757: 487b: Homo sapiens: RNA-Seq'.

SRA-BLAST does not currently support reference compressed SRA datasets, so it is generally advised that you add the condition 'NOT sra_nucore_alignment[Filter]' (as in the above example) to your queries to remove these datasets from the search results. Attempting to send incompatible datasets to BLAST will result in an error like the following:

Failed to convert SRX SRX079566 to SRA runs
Invalid SRX accession(s): SRX079566

If you believe that the data you are attempting to search against should be BLAST-able, but are not, please email sra@ncbi.nlm.nih.gov for assistance and advice. After successfully sending accessions to SRA BLAST, you are then able to input your sequence(s) of interest and perform the search.

NCBI/BLAST/BLASTN suite interface. The 'Enter Query Sequence' section is active, showing a text input field for the query sequence. The 'Choose Search Set' section shows a list of SRA Experiment sets (SRX) with a search bar. The 'Program Selection' section shows options for optimizing the search (Highly similar sequences, More dissimilar sequences, Somewhat similar sequences). The 'BLAST' button is highlighted.

Building a search list in SRA BLAST directly

SRA BLAST can be accessed directly. You will then need to provide SRA Experiment (SRX, DRX, or ERX) accessions or use the autocomplete feature to help refine your search. You may enter 1 Experiment accession per

line in the search list. The '+' button can then be used to add additional sequencing libraries to the search space. Note that a running tally of the number of sequences is presented above the list of accessions. Again, there is currently a limit of approximately 2 billion sequences per individual SRA BLAST query.

NCBI/ BLAST/ blastn suite **Sequence Read Archive Nucleotide BLAST** [Status of the NCBI Sequence Read Archive \(SRA\)](#)

blastn

Enter Query Sequence

BLASTN programs search SRA databases using a nucleotide query.

Enter accession number(s), gi(s), or FASTA sequence(s) [Clear](#) [Query subrange](#)

From

To

Or, upload file No file chosen

Job Title

Enter a descriptive title for your BLAST search

Choose Search Set

SRA Experiment set (SRX) **Sequences: 85,017,716**

Program Selection

Optimize for

BLAST

[Algorithm parameters](#)

SRX47

SRX470016 Amborella trichopoda bisulphite sequencing (Amborella trichopoda taxid:...

SRX470033 GSM1327148: H3K9me3_ChIPSeq (Ctrl); Mus musculus; ChIP-Seq (Mu...

SRX470034 GSM1327149: input DNA (Ctrl); Mus musculus; ChIP-Seq (Mus musculu...

SRX470035 GSM1327150: H3K9me3_ChIPSeq (Kd); Mus musculus; ChIP-Seq (Mus...

SRX470036 GSM1327151: input DNA (Kd); Mus musculus; ChIP-Seq (Mus musculu...

SRX470047 Whole genome shotgun sequencing of Salmonella enterica subsp. salam...

SRX470048 Whole genome shotgun sequencing of Salmonella enterica subsp. enteri...

SRX470049 Whole genome shotgun sequencing of Salmonella enterica subsp. enteri...

SRX470050 Whole genome shotgun sequencing of Salmonella enterica subsp. enteri...

SRX470051 Whole genome shotgun sequencing of Salmonella enterica subsp. enteri...

SRX470052 Whole genome shotgun sequencing of Salmonella enterica subsp. enteri...

SRX470053 Whole genome shotgun sequencing of Salmonella enterica subsp. salam...

0 top suggestions will be shown.

highlighted in yellow and marked with a sign

Direct downloading of fasta and fastq format

The SRA provides a [tool](#) that can be used to download data directly in fasta or fastq format. You must provide one or more SRA Experiment (SRX, DRX, or ERX) accessions in a comma-separated list. The same filtering inputs available in the Run Browser (described above) are available here to restrict the number of returned reads. Certain reads can also be clipped to remove low quality data from the download. If more than one Run accession in the list is checked, all data will be downloaded into a single fasta or fastq file, rather than per-accession files. Note that the output format of this tool is pre-defined and cannot be adjusted at the time of download. Users with specific formatting needs (e.g., for downstream analysis) are encouraged to use the [SRA Toolkit](#) to download and convert the data (described above).

NCBI [Site map](#) [All databases](#) [Search](#)

Sequence Read Archive

[Main](#) [Browse](#) [Search](#) [Download](#) [Submit](#) [Documentation](#) [Software](#) [Trace Archive](#) [Trace Assembly](#) [Trace Home](#) [Trace BLAST](#)

FASTA/FASTQ [Reads](#) [Analyses](#) [Reports](#) [References](#)

Download for Experiment SRX079566

Accession	# of bases	# of spots
<input type="checkbox"/> select all		total filtered
<input checked="" type="checkbox"/> SRR292241	699.9 M	9.7 M
<input checked="" type="checkbox"/> SRR390728	516.9 M	7.2 M

Filter

Search: [Apply](#)

[What can the filter be applied to?](#)

Download Format

☐ filtered ☐ clipped ☒ FASTA ☐ FASTQ [Download](#)

Downloading metadata associated with SRA data files

SRA data files do not contain any information about the metadata (sample information, etc.) linked to the data themselves. The SRA provides a few tools to allow downloading of metadata in batch. Note that these tools differ

from the Entrez [Experiment](#), [BioSample](#), and [BioProject](#) reports for a given dataset and may not contain all relevant metadata.

Viewing and downloading tabular metadata with the SRA Run Selector

The [SRA Run Selector](#) can be used to view metadata from one or more projects (SRA Study accessions – SRP, DRP, or ERP) entered into the field at the top of the page. The Run Selector provides a table view of library preparation and sample attribute metadata. The table can be filtered by sample attribute(s), accessions, etc. The “Get Metadata” button can be used to download a table (.txt, tab-delimited) of all or selected metadata.

NCBI

SRA RUN Selector

Change Study

SRP001599

Change

Common attributes:

SRA Study

BioProject

analyte_type

gap_accession

is_tumor

study_name

Assay Type

Center Name

Platform

Conse

SRP001599

phs000235

RNA

phs000235

1

NCI Cancer Genome Characterization Initiative (CGCI)

RNA-Seq

BCCAGSC

ILLUMINA

public

Runs

BioSamples

MBytes

MBases

PermaLink

Total

20

10

13,363

22,364

Get Metadata

Filtered

20

10

13,363

22,364

+

-

x

Selected

0

0

0

0

Show selected

x

Get Metadata

	Run	BioSample	Sample Name	SRA Sample	DSMZ	body_site	cell_line	sex	study_subject_id	Library Name	MBases	MBytes
<input type="checkbox"/>	SRR292240	SAMN00630373	HS0685	SRS212580	ACC 47	pleural effusion fluid	DOHH-2	male	DOHH-2	HS0685	395	239
<input type="checkbox"/>	SRR292241	SAMN00630374	HS0798	SRS212581	ACC 539	peritoneal cavity fluid	DB	male	DB	HS0798	667	958
<input type="checkbox"/>	SRR292242	SAMN00630375	HS0841	SRS212582	ACC 32	pleural effusion fluid	Karpas 422	female	Karpas 422	HS0841	718	447
<input type="checkbox"/>	SRR292243	SAMN00630376	HS0842	SRS212583	ACC 583	lymph node	NU-DHL-1	male	NU-DHL-1	HS0842	681	526
<input type="checkbox"/>	SRR292244	SAMN00630377	HS0900	SRS212584	ACC 572	peritoneal cavity fluid	SU-DHL-6	male	SU-DHL-6	HS0900	766	691
<input type="checkbox"/>	SRR292245	SAMN00630378	HS0901	SRS212585	ACC 575	pleural effusion fluid	WSU-DLCL2	male	WSU-DLCL2	HS0901	896	1,053
<input type="checkbox"/>	SRR292246	SAMN00630379	HS1163	SRS212586	ACC 722	bone marrow	OCI-LY1	male	OCI-LY1	HS1163	1,775	1,386
<input type="checkbox"/>	SRR292247	SAMN00630380	HS1182	SRS212587	ACC 688	peripheral blood	OCI-LY7	male	OCI-LY7	HS1182	2,140	1,603
<input type="checkbox"/>	SRR292248	SAMN00630381	HS1183	SRS212588	ACC-528	bone marrow	OCI-Ly19-R1	female	OCI-Ly19-R1	HS1183	1,619	1,265
<input type="checkbox"/>	SRR292249	SAMN00630382	HS2047	SRS212589	ACC 579	cerebrospinal fluid	NU-DUL-1	male	NU-DUL-1	HS2047	3,462	1,454
<input type="checkbox"/>	SRR390726	SAMN00630382	HS2047	SRS212589	ACC 579	cerebrospinal fluid	NU-DUL-1	male	NU-DUL-1	HS2047	3,394	821
<input type="checkbox"/>	SRR390727	SAMN00630373	HS0685	SRS212580	ACC 47	pleural effusion fluid	DOHH-2	male	DOHH-2	HS0685	342	137
<input type="checkbox"/>	SRR390728	SAMN00630374	HS0798	SRS212581	ACC 539	peritoneal cavity fluid	DB	male	DB	HS0798	492	184
<input type="checkbox"/>	SRR390729	SAMN00630375	HS0841	SRS212582	ACC 32	pleural effusion fluid	Karpas 422	female	Karpas 422	HS0841	602	143
<input type="checkbox"/>	SRR390730	SAMN00630376	HS0842	SRS212583	ACC 583	lymph node	NU-DHL-1	male	NU-DHL-1	HS0842	442	180
<input type="checkbox"/>	SRR390731	SAMN00630377	HS0900	SRS212584	ACC 572	peritoneal cavity fluid	SU-DHL-6	male	SU-DHL-6	HS0900	501	185
<input type="checkbox"/>	SRR390732	SAMN00630378	HS0901	SRS212585	ACC 575	pleural effusion fluid	WSU-DLCL2	male	WSU-DLCL2	HS0901	557	238
<input type="checkbox"/>	SRR390733	SAMN00630379	HS1163	SRS212586	ACC 722	bone marrow	OCI-LY1	male	OCI-LY1	HS1163	1,021	634
<input type="checkbox"/>	SRR390734	SAMN00630380	HS1182	SRS212587	ACC 688	peripheral blood	OCI-LY7	male	OCI-LY7	HS1182	917	613
<input type="checkbox"/>	SRR390735	SAMN00630381	HS1183	SRS212588	ACC-528	bone marrow	OCI-Ly19-R1	female	OCI-Ly19-R1	HS1183	977	606

Command line access to metadata with the SRA Run Info CGI

Users can access the SRA Run Info CGI either through a browser or using a command line tool like wget.

```
wget -O <file_name.csv> 'http://trace.ncbi.nlm.nih.gov/Traces/sra/sra.cgi?
save=efetch&db=sra&rettype=runinfo&term=<query>'
```

As a parallel to the above example in the Run Selector,

```
wget -O ./SRP001599_info.csv 'http://trace.ncbi.nlm.nih.gov/Traces/sra/sra.cgi?
save=efetch&db=sra&rettype=runinfo&term= SRP001599'
```

Will return essentially the same information. Note that the CGI returns data in a comma-separated value (.csv) format, rather than the tab-delimited format of the Run Selector. The last component, <query>, can contain any set of Entrez parameters. Users may refine a search using Entrez and then copy over the search terms to a script for batch downloading. As an example, the search string

"Homo sapiens"[Organism] AND "cancer"[All Fields] AND "cluster_public"[prop] AND "strategy wgs"[Properties]

Will return [these results](#) in an Entrez search of the SRA. The equivalent Run Info CGI search would be

```
wget -O ./query_results.csv 'http://trace.ncbi.nlm.nih.gov/Traces/sra/sra.cgi?
save=efetch&db=sra&rettype=runinfo&term="Homo sapiens"[Organism] AND "cancer"[All Fields]
AND "cluster_public"[prop] AND "strategy wgs"[Properties]'
```

Note that Entrez groups by Experiment accession, but that the CGI does not. It is, therefore, to be expected that the Run Info CGI will return a longer list of results than Entrez, but will still contain the same datasets.

Aspera Transfer Guide

Created: May 11, 2009; Updated: April 16, 2014.

Notice

Reference herein to any specific commercial products, process, or service by trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government, and shall not be used for advertising or product endorsement purposes.

Overview

This document provides instructions on the use and installation of Aspera Connect for high throughput file transfer with NCBI. As the sizes of the datasets have increased, we have found that the traditional methods of *ftp* or *http* do not have the performance characteristics needed to support this load of data.

Requirements for large scale data transfer over the internet include high bandwidth, auto checksum, recursive copy, and security based on strong keys. NCBI has chosen to use a product from Aspera, Inc (Emeryville, CA) because of improved data transfer characteristics. FTP and HTTP access will continue to be available and are the default options for users without Aspera installed. Instructions are provided below for investigators to use this data transfer technology. NCBI also is open to using additional products with the appropriate performance characteristics.

Scope

This document is intended for users transferring large data files to and from NCBI. It applies to the Sequence Read Archive (SRA), dbGaP, and other archives where aspera download is enabled.

Aspera

Aspera Connect

Aspera Connect is software that allows download and upload via a web plugin for popular browsers on machines running Linux, Windows, and Macintosh. The software also includes a command line tool (ascp) that allows scripted data transfer. The software client is free for users exchanging data with NCBI.

Download and install Aspera Connect software from: <http://downloads.asperasoft.com/connect2/>

The website's download button will default to the detected operating system of the user's computer. To download for a different OS, click the link to 'See all installers'.

Please note the Requirements and consult with your network administrator to ensure transfers with aspera will not be blocked.

Aspera can be installed for individual users. However users of shared machine may want to have the software installed for all users by a system administrator.

The fasp Protocol

The FASP protocol from Aspera (www.asperasoft.com) uses UDP, eliminating the latency issues seen with TCP, and provides bandwidth up to 5 gigabit per second (Gbps) to transfer data. It has a restart capability if data transfer is interrupted midstream and is well behaved, so if there is other data traffic on your network

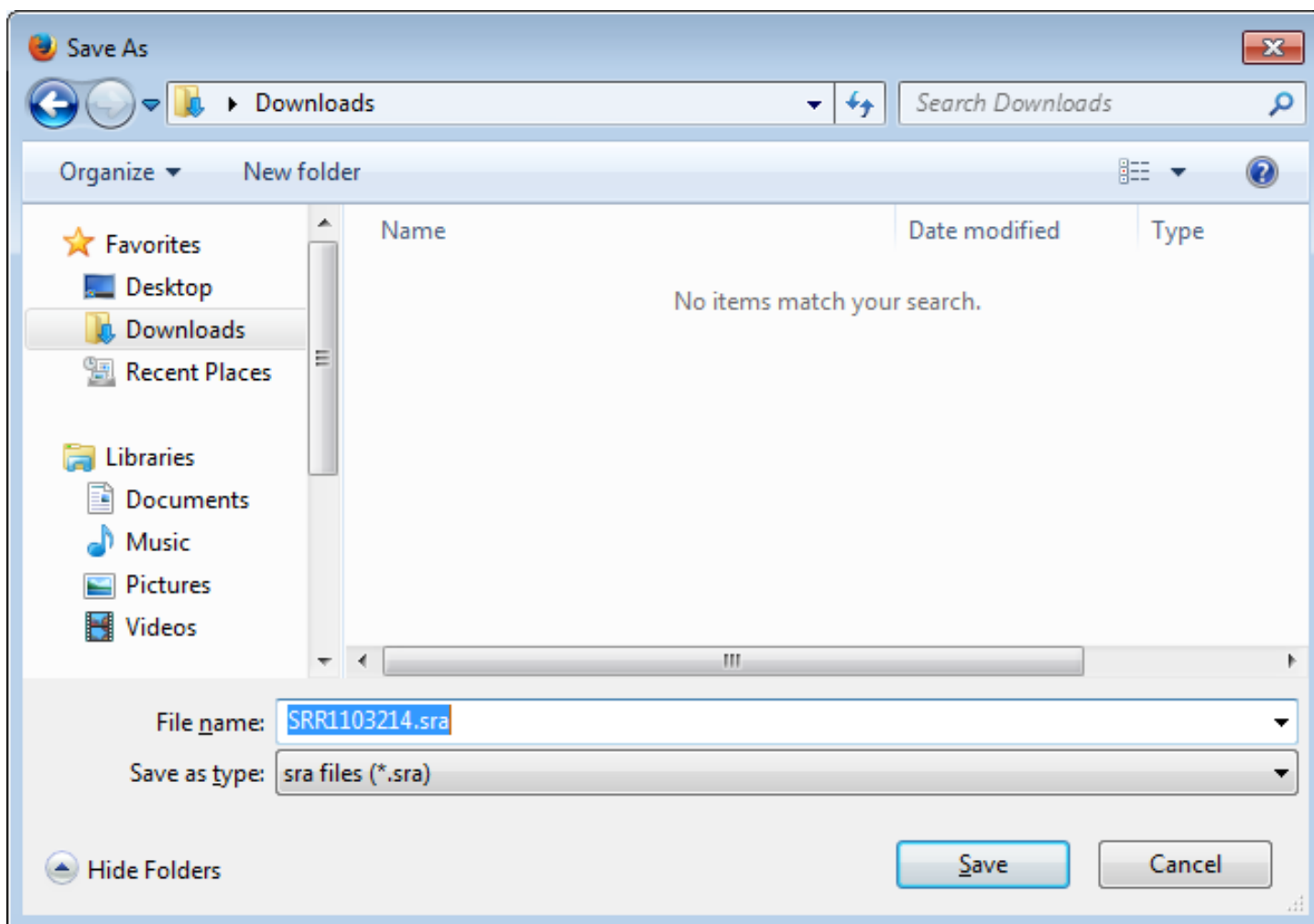
connections, it will back off in order to avoid starving other protocols. We have seen effective throughput up to 800 megabits per second (Mbps) to a single site.

Downloading Data with Aspera Connect Browser Plugin

Once the plugin has been installed in your browser, you may download files or entire directories from NCBI using Aspera. Example: In your browser window, go to

<http://www.ncbi.nlm.nih.gov/public/?ftp/sra/sra-instant/reads/ByRun/sra/SRR/SRR292/SRR292241>

Click 'SRR292241.sra' to begin saving the data. You will be prompted to select where the file is to be saved. For example:



You can download full directories or a single file at a time. The Aspera Connect plugin works with Chrome, Internet Explorer (IE), Safari, and FireFox web browsers. In some cases Aspera Connect may create a popup window to get a confirmation for file transfer and this popup window can be hidden behind your current web browser.

Using ascp to Download by Command Line

The command line program *ascp* is a utility delivered along with the Aspera Connect product.

```
ascp -i <asperaweb_id_dsa.openssh with path> -k1 -Tr -l100m  
anonftp@ftp.ncbi.nlm.nih.gov:/<files to transfer> <local destination>
```

- `-i <asperaweb_id_dsa.openssh with path>` = fully qualified path & file name where

this public key file is located. This file is part of Aspera Connect distribution and is usually located in the 'etc' subdirectory.

- `-T` to disable encryption
- `-k 1` enables resume of partial transfers
- `-r` recursive copy
- `-l` (maximum bandwidth of request, try 100M and go up from there)

Experiment with transfers starting at 100 Mbps and working up to 400 Mbps. Select the bandwidth setting that gives good performance with unattended operation.

- `<files(s) to transfer>` = names of files to transfer (including path)
- `<local destination path>` = location to store the downloaded data

Windows Executable Location

The *ascp* program for Microsoft Windows is located by default in "*C:\Program Files\Aspera\Aspera Connect\bin\ascp.exe*"

OS X Executable Location

The *ascp* Mac program location is */Applications/Aspera Connect.app/Contents/Resources/ascp*

Linux Executable Location

The *ascp* Linux program location is */opt/aspera/bin/ascp*

Additional information is available at the Aspera Web site: <http://downloads.asperasoft.com/documentation/>

Using ascp to Upload by Command Line

In order to use the Aspera upload service you will need to use a **private** SSH key, individual users can contact us at sra@ncbi.nlm.nih.gov to request an Aspera private key.

Upload Command

```
ascp -i <private key file> -T -l 100m <file(s) to transfer>
asp-****@upload.ncbi.nlm.nih.gov:<destination directory>
```

- `-i <private key file>` = fully qualified path & file name of the private SSH key
- `-T` to disable encryption
- `-k 1` enables resume of partial transfers
- `-l` (maximum bandwidth of request, try 100M and go up from there)

Experiment with transfers starting at 100 Mbps and working up to 400 Mbps. Select the bandwidth setting that gives good performance with unattended operation.

- `<files(s) to transfer>` = names of files to transfer (including path)
- `<destination directory>` = deposit location of the uploaded data (typically either 'test' or 'incoming')

For password protected private keys, it is possible to run *ascp* in an autonomous, unattended manner that does not require repeated login. The environmental variable `ASPERA_SCP_PASS` can be used to store the private key path for a scripted series of bulk uploads.

Key Pairs

SSH keys are used for establishing secure connections to remote computers.

Submitters using a dedicated center account can find instructions for generating a key pair or converting PuTTY format private keys to OpenSSH format in this guide.

<http://www.ncbi.nlm.nih.gov/books/NBK180157/>

Requirements

Firewall Requirements

Your local firewall must permit UDP data transfer in both directions on ports 33001-33009 for the following IP ranges:

130.14.*.*

165.112.*.*

The firewall must also allow ssh traffic outbound to NCBI.

Troubleshooting

Here are some example commands demonstrating a test download.

Mac OS X:

```
ascp -T -l640M -i "/Applications/Aspera Connect.app/Contents/Resources/
asperaweb_id_dsa.openssh" anonftp@ftp.ncbi.nlm.nih.gov:1GB /tmp/
```

Linux:

```
ascp -T -l640M -i /opt/aspera/etc/asperaweb_id_dsa.openssh
anonftp@ftp.ncbi.nlm.nih.gov:1GB /tmp/
```

MS Windows:

```
C:\TEMP>"C:\Program Files (x86)\Aspera\Aspera Connect\bin\ascp.exe" -T -l640M -
i "C:\Program Files (x86)\Aspera\Aspera Connect\etc\asperaweb_id_dsa.openssh" anon
ftp@ftp.ncbi.nlm.nih.gov:1GB C:\Temp\
```

For additional assistance, please contact the NCBI Help desk at info@ncbi.nlm.nih.gov

When you are about to contact the NCBI Help desk please provide them some basic information like operating system, version of aspera connect, type of disk storage used for transferring files and the type of network connection your organization has to the internet.

If you have a Linux or MacOS X operating system you may run these commands and show us their output:

```
curl -o /dev/null ftp://ftp.ncbi.nlm.nih.gov/1GB
curl -o /dev/null http://www.ncbi.nlm.nih.gov/staff/beloslyu/large.tar
traceroute ftp.ncbi.nlm.nih.gov
```

First two commands download a 1GB file from NCBI using ftp and http protocols, the content is dumped to /dev/null. The third command will let us see the latency in your internet connection and possible congestions on the way to NCBI.

Another possibility is to make some test downloads from Aspera's demo server, for Linux the command line is:

```
env ASPERA_SCP_PASS=demoaspera ascp -L- -T -l100m aspera@demo.asperasoft.com:aspera-test-dir-large/1GB /tmp/
```

Aspera Connect is a commercial product and program specific support is available from the manufacturer at <http://asperasoft.com/support/>

The currently up-to-date documentation for ascp can be found at <http://downloads.asperasoft.com/en/documentation/8>