

MCB 5472

Lecture #1 – NCBI and public data

Jan 27/14

Pretest question #1

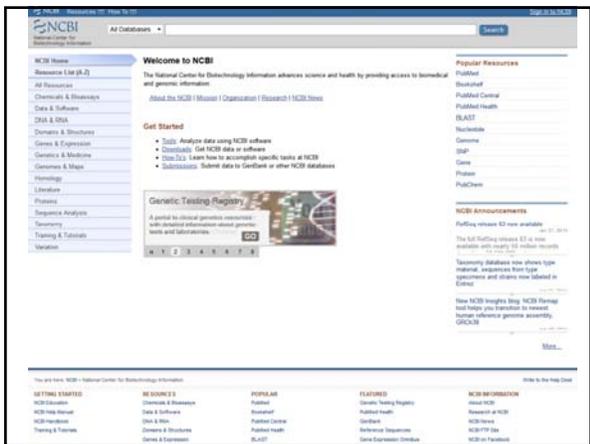
- Define “Evolution by Natural Selection” in one or two sentences
 - Variation exists in a population
 - Variation is inherited
 - Individuals reproduce more than the environment can support, therefore competition occurs
 - Differential reproduction is driven by heritable variation allowing more efficient resource use

Pretest question #2

- Can two sequences that are aligned reasonably well over their whole lengths be 66% homologous?
 - Sequences either are or are not homologous.
 - Homologous sequences can vary in their similarity

NCBI

- www.ncbi.nlm.nih.gov
- National Center for Biotechnology Information
- Run by the US National Institute of Health National Library of Medicine
- Since 1988 has been the designated US repository for data relating to biological research
- Also runs its own research program in computational biology



Comprised of 42 databases containing together >900 million records

Table 1. The NCBI database as of 1 September 2013

Database	Records	Section within this article	Data source
NCBI Web Site	21,929	Introduction	N
PubMed	23,402,796	Literature	C
PMID	13,626,362	Literature	D, C
NCBI e-library	1,407,000	Literature	C, N
Medline	1,011,719	Literature	N
Bookshelf	112,216	Literature	C, N
Genbank ^a	1,113,761	Genomes and RNA	C, N
RefSeq ^a	101,797,766	Genomes and RNA	D, GenBank, C, N
EMBL ^a	14,917,096	Genomes and RNA	D, GenBank, C, N
Gen	1,000,000	Genomes and RNA	C, N
Protein	2,088,017	Genomes and RNA	D
PubChem	69,661,474	Genomes and RNA	D, PubRank
Protein Families ^a	82,821	Proteins	N
NCBI Pathway ^a	76,797,791	Genes and expression	D
Pubx	11,367,498	Genes and expression	D
Gene ^a	14,142,886	Genes and expression	C, N
UniProt ^a	6,467,083	Genes and expression	N
NCBI Taxonomy ^a	1,044,144	Genes and expression	N
Biocompare ^a	122,217	Genes and expression	C, N
PhosphoSitePlus ^a	111,548	Genes and expression	N
ChEMBL ^a	13,113,797	Chemicals	D, N
ChEMBL	68,118	Chemicals	D
PubMeds ^a	17,707	Chemicals	C, N
Chemical	19,973	Chemicals	C, N
PubMeds ^a	19,973	Chemicals	D
SNP	100,100,000	Genetics and medicine	D, dbSNP, N
dbVar ^a	1,000,000	Genetics and medicine	D
MedGen ^a	100,000	Genetics and medicine	C, N
dbCSP ^a	154,971	Genetics and medicine	D
dbCSP	69,061	Genetics and medicine	D, N
PubMed Health	41,267	Genetics and medicine	C
dbCSP	21,211	Genetics and medicine	D
dbCSP	20,056	Genetics and medicine	C
PubChem Substances ^a	119,113,846	Chemicals and bioassays	D
PubChem Compounds ^a	47,717,096	Chemicals and bioassays	D
PubChem Bioassays ^a	171,429	Chemicals and bioassays	D
Structure ^a	62,991	Chemicals and structures	C, N
EMBL	68,014	Chemicals and structures	C, N

^aIndicates that the data in this resource are available by FTP.
^b dbSNP accession; C, eLibrary/Entrez; N, National Center for Biotechnology Information.

NCBI Resource Coordinators Nucleic Acids Res. 2014 42: D7-D17

GenBank

- Central database for nucleotide sequences
- Synced daily with alternative independent databases
 - European Nucleotide Archive (ENA)
 - DNA Data Bank of Japan (DDBJ)
- Linked to other NCBI resources via Entrez system

GenBank sequence types

- WGS (whole genome sequencing)
- TSA (transcriptome shotgun assembly)
- 12 standard nucleotide types
 - Split by taxon + ENV for environmental sequences
- PAT (patent sequences)
- 6 high-throughput types
 - EST (expressed sequence tag)
 - GSS (genome survey sequence)
 - HTC (high-throughput cDNA)
 - HTG (high-throughput genomic)
 - STS (sequence tagged sites)

NCBI raw nucleotide data

- SRA – Sequence Read Archive: raw next-generation sequencing data, typically underpins a GenBank genome assembly
- Trace Archive: raw Sanger sequencing data

Table 1. Growth of GenBank divisions (nucleotide base pairs)

Division	Description	Release 197 (8/2013)	Annual increase (%) ^a
WGS	Whole-genome shotgun data	500 420 412 665	62.4
TSA	Transcriptome shotgun data	8 633 123 935	49.9
PHG	Phages	119 812 712	42.5
VRL	Viruses	1 757 202 472	22.9
BCT	Bacteria	10 281 048 518	21.8
ENV	Environmental samples	3 743 277 434	10.9
INV	Invertebrates	2 737 140 646	9.8
PAT	Patented sequences	13 290 161 247	9.7
PLN	Plants	5 963 882 822	8.8
GSS	Genome survey sequences	23 726 384 753	8.1
VRT	Other vertebrates	3 068 956 026	6.3
MAM	Other mammals	911 342 025	5.6
HTG	High-throughput genomic	25 184 819 955	3.4
HTC	High-throughput cDNA	656 196 063	2.7
UNA	Unannotated	130 510	2.1
EST	Expressed sequence tags	41 665 629 009	1.9
PRI	Primates	6 425 093 034	1.7
SYN	Synthetic	941 078 074	1.4
ROD	Rodents	4 451 315 297	0.4
STS	Sequence tagged sites	636 326 479	0.0
TOTAL	All GenBank sequences	654 613 333 676	45.1

^aMeasured relative to Release 191 (8/2012).

Benson et al. Nucl. Acids Res. 2014 42: D32-D37

RefSeq

- RefSeq is a curated collection of annotated sequences
- Typically, raw nucleotide sequences come from GenBank
- Annotations are typically from NCBI, but sometimes from collaborations with the broader community
- Some quality restrictions for RefSeq annotating GenBank data

Discuss

1. What makes a database? Why use a database, e.g., vs. a spreadsheet
2. GenBank is a repository, but RefSeq is not. Why is this significant?

Discuss

3. What data types of genomic data are there?

4. How are these data types related

.gbk files

- NCBI's standard format of annotating nucleotide sequences

```
LOCUS       M_12094          1472 bp     DNA     linear     BP     21-MAY-2011
DEFINITION  Escherichia coli str. F-12 substr. M01455 strain F-12 145 ribosomal
            RNA, complete sequence.
ACCESSION   M_12094
VERSION    M_12094.1  01/07/1997
BUILD      RefSeqProject PRNA33179
KEYWORDS   RefSeq
SOURCE     Escherichia coli str. F-12 substr. M01455
ORGANISM   Escherichia coli str. F-12 substr. M01455
            Bacteria; Proteobacteria; Gammaproteobacteria; Enterobacteriales;
            Enterobacteriaceae; Escherichia.
REFERENCE  1 (base 1 to 1472)
            Riley,M., Amst., Arnold,M.B., Berlyn,M.K., Blattner,F.R.,
            Chaudhuri,K.K., Glanzer,J.D., Horiuchi,T., Kessler,I.M., Kowop,T.,
            Muz,W., Perna,M.T., Plunkett,G. III, Rudi,K.E., Serres,M.H.,
            Thomas,S.H., Thomson,M.B., Whittam,D. and Wanner,B.L.
            Escherichia coli F-12: a cooperatively developed annotation
            mapshot-2005
JOURNAL    Nucleic Acids Res. 34 (1), 1-9 (2006)
PUBMED    1437929
REMARK     Publication Status: Online-Only
AUTHORS   Blattner,F.R. and Plunkett,G. III.
TITLE      Direct Submission
JOURNAL    Submitted (16-JAN-1997) Laboratory of Genetics, University of
            Wisconsin, 4305 Henry Hall, Madison, WI 53706-1580, USA
COMMENT   REVISED REFSEQ: This record has been curated by NCBI staff. The
            reference sequence is identical to 000994|22371-22312.
            This record has been curated by collaboration with an international
            consortium of ribosomal RNA databases.
            COMPLETES: full length.
PRIMARY   REFSEQ_SPAN      PRIMER1 IDENTIFIER PRIMER2 SPAN      COMP
            1-1542              000994.2      22371-22312
FEATURES             Location/Qualifiers
             source          1..1542
                        /organism="Escherichia coli str. F-12 substr. M01455"
                        /mol_type="rRNA"
                        /rname="r12"
                        /rdb_xref="rM01455"
                        /db_xref="taxon:511418"
             rRNA           1..1542
                        /product="14S ribosomal RNA"
ORIGIN          1 aaattgaga gtttgatct gtttcattt gaattgttg gttgagctta aactctgaa
            41 gttgagcttg aataggaga agttctgctt ttgtgtgag agttgagag gttctgctaa
            ...
            1461 agttacttc tttggagga gttactcttc ttgtatgag tgaattggt gaattctaa
            1501 caattgaaq ttgagggag ttgtgtgttg atctactct ta
            //
```

Accession: NCBI's identifier for this sequence record

Version: NCBI's identifier for this sequence record version

GI: NCBI's identifier for this unique document

These features are the only ones that are "stable", i.e., absolute identifiers

```
LOCUS       M_12094          1472 bp     DNA     linear     BP     21-MAY-2011
DEFINITION  Escherichia coli str. F-12 substr. M01455 strain F-12 145 ribosomal
            RNA, complete sequence.
ACCESSION   M_12094
VERSION    M_12094.1  01/07/1997
BUILD      RefSeqProject PRNA33179
KEYWORDS   RefSeq
SOURCE     Escherichia coli str. F-12 substr. M01455
ORGANISM   Escherichia coli str. F-12 substr. M01455
            Bacteria; Proteobacteria; Gammaproteobacteria; Enterobacteriales;
            Enterobacteriaceae; Escherichia.
REFERENCE  1 (base 1 to 1472)
            Riley,M., Amst., Arnold,M.B., Berlyn,M.K., Blattner,F.R.,
            Chaudhuri,K.K., Glanzer,J.D., Horiuchi,T., Kessler,I.M., Kowop,T.,
            Muz,W., Perna,M.T., Plunkett,G. III, Rudi,K.E., Serres,M.H.,
            Thomas,S.H., Thomson,M.B., Whittam,D. and Wanner,B.L.
            Escherichia coli F-12: a cooperatively developed annotation
            mapshot-2005
JOURNAL    Nucleic Acids Res. 34 (1), 1-9 (2006)
PUBMED    1437929
REMARK     Publication Status: Online-Only
AUTHORS   Blattner,F.R. and Plunkett,G. III.
TITLE      Direct Submission
JOURNAL    Submitted (16-JAN-1997) Laboratory of Genetics, University of
            Wisconsin, 4305 Henry Hall, Madison, WI 53706-1580, USA
COMMENT   REVISED REFSEQ: This record has been curated by NCBI staff. The
            reference sequence is identical to 000994|22371-22312.
            This record has been curated by collaboration with an international
            consortium of ribosomal RNA databases.
            COMPLETES: full length.
PRIMARY   REFSEQ_SPAN      PRIMER1 IDENTIFIER PRIMER2 SPAN      COMP
            1-1542              000994.2      22371-22312
FEATURES             Location/Qualifiers
             source          1..1542
                        /organism="Escherichia coli str. F-12 substr. M01455"
                        /mol_type="rRNA"
                        /rname="r12"
                        /rdb_xref="rM01455"
                        /db_xref="taxon:511418"
             rRNA           1..1542
                        /product="14S ribosomal RNA"
ORIGIN          1 aaattgaga gtttgatct gtttcattt gaattgttg gttgagctta aactctgaa
            41 gttgagcttg aataggaga agttctgctt ttgtgtgag agttgagag gttctgctaa
            ...
            1461 agttacttc tttggagga gttactcttc ttgtatgag tgaattggt gaattctaa
            1501 caattgaaq ttgagggag ttgtgtgttg atctactct ta
            //
```

Locus: A name for this stretch of DNA

Length of DNA

Molecule type

GenBank division

Date of last modification

```
LOCUS       M_12094          1472 bp     DNA     linear     BP     21-MAY-2011
DEFINITION  Escherichia coli str. F-12 substr. M01455 strain F-12 145 ribosomal
            RNA, complete sequence.
ACCESSION   M_12094
VERSION    M_12094.1  01/07/1997
BUILD      RefSeqProject PRNA33179
KEYWORDS   RefSeq
SOURCE     Escherichia coli str. F-12 substr. M01455
ORGANISM   Escherichia coli str. F-12 substr. M01455
            Bacteria; Proteobacteria; Gammaproteobacteria; Enterobacteriales;
            Enterobacteriaceae; Escherichia.
REFERENCE  1 (base 1 to 1472)
            Riley,M., Amst., Arnold,M.B., Berlyn,M.K., Blattner,F.R.,
            Chaudhuri,K.K., Glanzer,J.D., Horiuchi,T., Kessler,I.M., Kowop,T.,
            Muz,W., Perna,M.T., Plunkett,G. III, Rudi,K.E., Serres,M.H.,
            Thomas,S.H., Thomson,M.B., Whittam,D. and Wanner,B.L.
            Escherichia coli F-12: a cooperatively developed annotation
            mapshot-2005
JOURNAL    Nucleic Acids Res. 34 (1), 1-9 (2006)
PUBMED    1437929
REMARK     Publication Status: Online-Only
AUTHORS   Blattner,F.R. and Plunkett,G. III.
TITLE      Direct Submission
JOURNAL    Submitted (16-JAN-1997) Laboratory of Genetics, University of
            Wisconsin, 4305 Henry Hall, Madison, WI 53706-1580, USA
COMMENT   REVISED REFSEQ: This record has been curated by NCBI staff. The
            reference sequence is identical to 000994|22371-22312.
            This record has been curated by collaboration with an international
            consortium of ribosomal RNA databases.
            COMPLETES: full length.
PRIMARY   REFSEQ_SPAN      PRIMER1 IDENTIFIER PRIMER2 SPAN      COMP
            1-1542              000994.2      22371-22312
FEATURES             Location/Qualifiers
             source          1..1542
                        /organism="Escherichia coli str. F-12 substr. M01455"
                        /mol_type="rRNA"
                        /rname="r12"
                        /rdb_xref="rM01455"
                        /db_xref="taxon:511418"
             rRNA           1..1542
                        /product="14S ribosomal RNA"
ORIGIN          1 aaattgaga gtttgatct gtttcattt gaattgttg gttgagctta aactctgaa
            41 gttgagcttg aataggaga agttctgctt ttgtgtgag agttgagag gttctgctaa
            ...
            1461 agttacttc tttggagga gttactcttc ttgtatgag tgaattggt gaattctaa
            1501 caattgaaq ttgagggag ttgtgtgttg atctactct ta
            //
```

Dblink: other NCBI databases related to this record

Source: the organism from which this sequence was determined

Organism: Taxonomic breakdown for the SOURCE organism, according to NCBI's Taxonomy DB

```
LOCUS       M_12094          1472 bp     DNA     linear     BP     21-MAY-2011
DEFINITION  Escherichia coli str. F-12 substr. M01455 strain F-12 145 ribosomal
            RNA, complete sequence.
ACCESSION   M_12094
VERSION    M_12094.1  01/07/1997
BUILD      RefSeqProject PRNA33179
KEYWORDS   RefSeq
SOURCE     Escherichia coli str. F-12 substr. M01455
ORGANISM   Escherichia coli str. F-12 substr. M01455
            Bacteria; Proteobacteria; Gammaproteobacteria; Enterobacteriales;
            Enterobacteriaceae; Escherichia.
REFERENCE  1 (base 1 to 1472)
            Riley,M., Amst., Arnold,M.B., Berlyn,M.K., Blattner,F.R.,
            Chaudhuri,K.K., Glanzer,J.D., Horiuchi,T., Kessler,I.M., Kowop,T.,
            Muz,W., Perna,M.T., Plunkett,G. III, Rudi,K.E., Serres,M.H.,
            Thomas,S.H., Thomson,M.B., Whittam,D. and Wanner,B.L.
            Escherichia coli F-12: a cooperatively developed annotation
            mapshot-2005
JOURNAL    Nucleic Acids Res. 34 (1), 1-9 (2006)
PUBMED    1437929
REMARK     Publication Status: Online-Only
AUTHORS   Blattner,F.R. and Plunkett,G. III.
TITLE      Direct Submission
JOURNAL    Submitted (16-JAN-1997) Laboratory of Genetics, University of
            Wisconsin, 4305 Henry Hall, Madison, WI 53706-1580, USA
COMMENT   REVISED REFSEQ: This record has been curated by NCBI staff. The
            reference sequence is identical to 000994|22371-22312.
            This record has been curated by collaboration with an international
            consortium of ribosomal RNA databases.
            COMPLETES: full length.
PRIMARY   REFSEQ_SPAN      PRIMER1 IDENTIFIER PRIMER2 SPAN      COMP
            1-1542              000994.2      22371-22312
FEATURES             Location/Qualifiers
             source          1..1542
                        /organism="Escherichia coli str. F-12 substr. M01455"
                        /mol_type="rRNA"
                        /rname="r12"
                        /rdb_xref="rM01455"
                        /db_xref="taxon:511418"
             rRNA           1..1542
                        /product="14S ribosomal RNA"
ORIGIN          1 aaattgaga gtttgatct gtttcattt gaattgttg gttgagctta aactctgaa
            41 gttgagcttg aataggaga agttctgctt ttgtgtgag agttgagag gttctgctaa
            ...
            1461 agttacttc tttggagga gttactcttc ttgtatgag tgaattggt gaattctaa
            1501 caattgaaq ttgagggag ttgtgtgttg atctactct ta
            //
```

Dblink: other NCBI databases related to this record

Source: the organism from which this sequence was determined

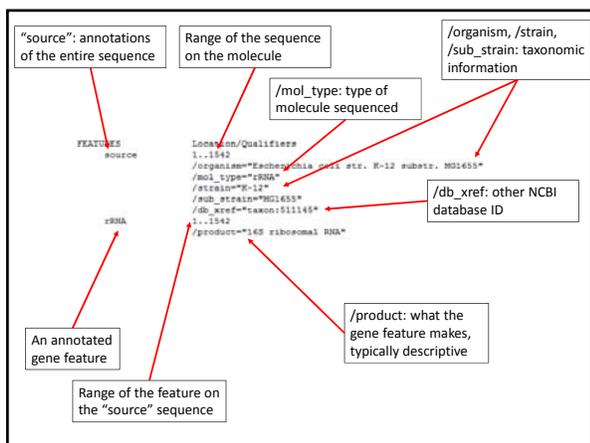
Organism: Taxonomic breakdown for the SOURCE organism, according to NCBI's Taxonomy DB

Reference: reference(s) relating to the sequence

Origin: The nucleotide sequence

"/": always marks the end of the file

```
LOCUS       M_12094          1472 bp     DNA     linear     BP     21-MAY-2011
DEFINITION  Escherichia coli str. F-12 substr. M01455 strain F-12 145 ribosomal
            RNA, complete sequence.
ACCESSION   M_12094
VERSION    M_12094.1  01/07/1997
BUILD      RefSeqProject PRNA33179
KEYWORDS   RefSeq
SOURCE     Escherichia coli str. F-12 substr. M01455
ORGANISM   Escherichia coli str. F-12 substr. M01455
            Bacteria; Proteobacteria; Gammaproteobacteria; Enterobacteriales;
            Enterobacteriaceae; Escherichia.
REFERENCE  1 (base 1 to 1472)
            Riley,M., Amst., Arnold,M.B., Berlyn,M.K., Blattner,F.R.,
            Chaudhuri,K.K., Glanzer,J.D., Horiuchi,T., Kessler,I.M., Kowop,T.,
            Muz,W., Perna,M.T., Plunkett,G. III, Rudi,K.E., Serres,M.H.,
            Thomas,S.H., Thomson,M.B., Whittam,D. and Wanner,B.L.
            Escherichia coli F-12: a cooperatively developed annotation
            mapshot-2005
JOURNAL    Nucleic Acids Res. 34 (1), 1-9 (2006)
PUBMED    1437929
REMARK     Publication Status: Online-Only
AUTHORS   Blattner,F.R. and Plunkett,G. III.
TITLE      Direct Submission
JOURNAL    Submitted (16-JAN-1997) Laboratory of Genetics, University of
            Wisconsin, 4305 Henry Hall, Madison, WI 53706-1580, USA
COMMENT   REVISED REFSEQ: This record has been curated by NCBI staff. The
            reference sequence is identical to 000994|22371-22312.
            This record has been curated by collaboration with an international
            consortium of ribosomal RNA databases.
            COMPLETES: full length.
PRIMARY   REFSEQ_SPAN      PRIMER1 IDENTIFIER PRIMER2 SPAN      COMP
            1-1542              000994.2      22371-22312
FEATURES             Location/Qualifiers
             source          1..1542
                        /organism="Escherichia coli str. F-12 substr. M01455"
                        /mol_type="rRNA"
                        /rname="r12"
                        /rdb_xref="rM01455"
                        /db_xref="taxon:511418"
             rRNA           1..1542
                        /product="14S ribosomal RNA"
ORIGIN          1 aaattgaga gtttgatct gtttcattt gaattgttg gttgagctta aactctgaa
            41 gttgagcttg aataggaga agttctgctt ttgtgtgag agttgagag gttctgctaa
            ...
            1461 agttacttc tttggagga gttactcttc ttgtatgag tgaattggt gaattctaa
            1501 caattgaaq ttgagggag ttgtgtgttg atctactct ta
            //
```



LOCUS	NT_077402	257719 bp	DNA	linear	CON 13-AUG-2013
DEFINITION	Homo sapiens chromosome 1 genomic contig, GRCh37.p13 Primary Assembly				
ACCESSION	NT_077402 GDS_000125215 NT_077911				
VERSION	NT_077402.2 GI:224514618				
DBLINK	BioProject: PRJNA1468				
KEYWORDS	RefSeq				
SOURCE	Homo sapiens (human)				
ORGANISM	Homo sapiens				
REFERENCE	1 (bases 1 to 257719)				
AUTHORS	Gregory,S.D., Barlow,K.F., McLay,K.E., Kaul,R., Warburton,D., Dunham,A., Scott,C.E., Howe,K.L., Woodfine,K., Spencer,C.C., Jones,M.C., Gillson,C., Searle,S., Zhou,Y., Kokocinski,F., McDonald,S., Evans,K., Phillips,K., Akiyama,A., Cooper,R., Jones,C., Mall.R.E., Andrews,T.D., Lloyd,C., Ainscough,R., Almeida,J.P., Ambrose,K.D., Anderson,F., Andrew,R.W., Ashwell,R.F.I., Aubin,K., Babage,A.K., Baguley,C.L., Bailey,J., Baisley,H., Bethel,G., Bird,C.P., Bray-Allen,S., Brown,J.Y., Brown,A.J., Buckley,D., Burton,J., Bye,J., Carder,C., Chapman,J.C., Clark,S.Y., Clarke,G., Clem,C., Cudley,V., Collier,R.E., Corby,K., Coville,G.J., Davies,J., Deadman,R., Dunn,M., Barthrow,M., Ellington,A.G., Errington,H., Frankish,A., Frankland,J., French,L., Garner,P., Garnett,J., Gay,L., Ghori,M.R., Gibson,R., Gilby,L.M., Gillett,W., Glibert,R.J., Grafham,D.V., Griffiths,C., Griffiths-Jones,S., Grocock,R., Hammond,S., Harrison,E.S., Hart,E., Haugen,K., Heath,P.D., Holmes,S., Holt,K., Howden,P.J., Hunt,A.R., Hunt,S.K., Hunter,D., Isherwood,J., James,P., Johnson,C., Johnson,D., Joy,A., Kay,M., Kersey,J.K., Kibukawa,M., Kimberley,A.M., King,A., Knights,A.J., Lad,M., Laird,G., Lawlor,S., Leongamornrat,D.A., Lloyd,D.M., Loveland,J., Lovell,J., Lush,M.J., Lyne,R., Martin,S., Mashreghi-Mohammadi,M., Matthews,L., Matthews,N.S., McLaren,S., Milne,S., Mistry,S., Moore,S.P., Nicerson,T., O'Neil,C.N., Oliver,K., Palmieri,A., Palmer,S.A., Parker,A., Patel,D., Pearce,A.V., Peck,A.I., Pelan,S., Phelps,K., Phillimore,B.J., Plumb,R., Rajan,J., Raymond,C., Rouse,G., Sanchiriac,M., Schae,H.F., Sheridan,S., Shownkeen,S., Sims,S.				

Shane,C.D., Smith,M., Steward,C., Subramanian,S., Symons,M., Tracey,A., Tromans,A., Van Helmond,J., Wall,M., Wallis,J.M., White,S., Whitehead,S.L., Wilkinson,J.E., Willey,D.L., Williams,H., Wilming,L., Wray,P.W., Wu,Z., Coulson,A., Vaudin,W., Sulston,J.E., Durbin,S., Hubbard,T., Wooster,R., Dunham,I., Carter,N.P., McVean,G., Ross,M.T., Harrow,J., Olson,M.V., Beck,S., Rogers,J., Bentley,D.R., Banerjee,R., Bryant,S.P., Burford,D.C., Burrill,W.D., Clegg,S.M., Dhani,P., Dovey,O., Faulkner,L.W., Gribble,S.M., Langford,C.F., Pandian,R.D., Porter,K.M. and Prigmore,E.	
TITLE	The human genome project had a few authors!
JOURNAL	Nature 441 (7094): 315-321 (2004)
PubMed	16710414
REMARK	Erratum:[Nature, 2006 Oct 26;443(7141):1013. Banerjee, S [added]; Bryant, SP [added]; Burford, DC [added]; Burrill, WDS [added]; Clegg, SM [added]; Dhani, P [added]; Dovey, O [added]; Faulkner, LM [added]; Gribble, SM [added]; Langford, CF [added]; Pandian, RD [added]; Porter, KM [added]; Prigmore, E [added]]
REFERENCE	2 (bases 1 to 257719)
CONSTRM	International Human Genome Sequencing Consortium
TITLE	Finishing the euchromatic sequence of the human genome
JOURNAL	Nature 431 (7011): 931-945 (2004)
PubMed	15496913
REFERENCE	3 (bases 1 to 257719)
AUTHORS	Lander,E.S., Linton,L.M., Birren,B., Nusbaum,C., Zody,M.C., Baldwin,J., Devon,K., Dewar,K., Doyle,M., FitzHugh,M., Funke,R., Gage,D., Harris,K., Heaford,A., Howland,J., Kann,L., Lehoczky,J., LeVine,R., McEwan,P., McKernan,K., Meldrum,J., Mesirov,J.P., Miranda,C., Morris,M., Naylor,J., Raymond,C., Rosetti,M., Saitou,N., Sheridan,A., Sougier,C., Stange-Thomann,M., Stojanovic,M., Subramanian,A., Wyman,D., Rogers,J., Sulston,J., Ainscough,R., Beck,S., Bentley,D., Burton,J., Clee,C., Carter,N., Coulson,A., Deadman,R., Deloukas,P., Dunham,A., Dumbin,L., Durbin,S., French,L., Grafham,D., Gregory,S., Hubbard,T., Humphray,S., Hunt,A., Jones,M., Lloyd,C., Murray,A., Matthews,L., Mercer,S., Milne,S., Wallis,J.C., Mungall,A., Plumb,R., Ross,M., Showkeen,S., Sims,S., Waterston,R.H., Wilson,R.K., Hillier,L.M., McPherson,J.D., Marra,M.A., Mardis,E.R., Fulton,L.A., Chinwalla,A.T., Pigin,K.H., Ghisla,W.R., Chissoe,S.L., Wendl,M.C., Delehaunty,K.D., Miner,T.L., Delehaunty,A., Kramer,J.B., Cook,L.L., Fulton,B.S., Johnson,D.L., Mitx,P.J., Clifton,S.W., Hawkins,T.,

Brancomb,E., Reddi,P., Richardson,P., Henning,S., Slezak,T., Doggett,N., Cheng,J.F., Olsen,A., Lucas,S., Elkin,C., Uberbacher,E., Frazier,M., Gibbs,B.A., Muzny,D.M., Scherer,S.E., Souk,J.S., Sodergren,E.J., Worley,K.C., Zieve,C.M., Coriell,J.H., Metzker,M.L., Naylor,S.L., Kucherlapati,R.S., Nelson,D.L., Weinstock,G.M., Sakaki,Y., Fujiyama,A., Hattori,M., Yada,T., Toyoda,A., Itoh,T., Kawagoe,C., Watanabe,H., Totoki,Y., Taylor,T., Weissbach,J., Hellmuth,S., Burin,M., Artiguenave,F., Brentier,P., Bruns,T., Pelletier,E., Robert,C., Minckler,P., Smith,D.R., Doucette-Stamm,L., Rubinfeld,M., Weinstock,K., Lee,K.M., Dubois,J., Rosenblatt,J., Glavner,M., Baylata,D., Taudien,S., Rump,A., Yang,H., Yu,J., Wang,J., Huang,C., Gu,J., Hood,L., Rowen,L., Madan,A., Qin,S., Davis,R.M., Fedrizzi,M.A., Abola,A.P., Proctor,M.P., Myers,R.M., Schmutz,J., Bonham,C., Grimwood,J., Cox,D.R., Olson,M.V., Kaul,R., Raymond,C., Shimizu,M., Kawasaki,K., Minoshima,S., Evans,G.A., Athanasiou,M., Schultz,R., Roe,B.A., Chen,F., Fan,R., Roeber,J., Lavruchh,K., Reinhardt,R., McCombie,W.R., de la Bastide,M., Dedhia,M., Blocker,H., Hornischer,K., Nordliek,G., Agarwala,R., Aravind,L., Bailey,J.A., Bateman,A., Batzoglou,S., Birney,E., Bork,P., Brown,D.G., Burge,C.B., Cerutti,L., Chen,H.C., Church,D., Clamp,M., Copley,R.R., Doerks,T., Eddy,S.R., Eichler,E.E., Furey,T.S., Galagan,J., Gilbert,J.D., Harrow,C., Hayashizaki,Y., Haussler,D., Hermjakob,H., Hokamp,K., Jang,M., Johnson,L.S., Jones,T.A., Kasif,S., Kasprzyk,A., Kennedy,S., Kent,W.J., Kitts,P., Konin,K.V., Korfi,I., Kulp,D., Lancet,D., Lowe,T.M., Malyugant,A., Mikkelsen,T., Moran,J.V., Mulder,N., Pollar,V.J., Ponting,C.P., Schuler,G., Schultz,J., Slater,G., Smit,A.P., Stupka,K., Sutakowski,J., Thierry-Mieg,D., Thierry-Mieg,J., Wagner,L., Wallis,J., Wheeler,K., Williams,A., Wolf,Y.I., Wolfe,K.H., Yang,S.P., Yeh,R.F., Collins,F., Guyer,M.S., Peterson,J., Felsenfeld,A., Wetterstrand,P., Patrino,A., Morgan,M.J., de Jong,P., Catanese,J.J., Osoegawa,K., Shiyaya,R., Choi,S. and Chen,Y.	
CONSTRM	International Human Genome Sequencing Consortium
TITLE	Initial sequencing and analysis of the human genome
JOURNAL	Nature 409 (6822): 860-921 (2001)
PubMed	11237011
REMARK	Erratum:[Nature 2001 Aug 2;412(6846):565]

COMMENT	REFSEQ INFORMATION: The reference sequence is identical to GLO00001.1. On or before Jun 10, 2009 this sequence version replaced gi:29794400, gi:29794392. Assembly Name: GRCh37.p13 Primary Assembly The DNA sequence is composed of genomic sequence, primarily finished clones that were sequenced as part of the Human Genome Project. PCR products and WGS shotgun sequence have been added where necessary to fill gaps or correct errors. All such additions are manually curated by GRC staff. For more information see: http://genomereference.org.
##Genome-Annotation-Data-START##	Annotation Provider : NCBI
Annotation Status	: Full annotation
Annotation Version	: Homo sapiens Annotation Release 105
Annotation Pipeline	: NCBI eukaryotic genome annotation pipeline
Annotation Software Version	: 5.1
Annotation Method	: Best-placed RefSeq Genom
Features Annotated	: Gene; mRNA; CDS; ncRNA
##Genome-Annotation-Data-END##	
FEATURES	Location/Qualifiers
source	1..257719
	/organism="Homo sapiens"
	/mol_type="genomic DNA"
	/db_xref="taxon:9606"
	/chromosome="1"

gene	1874..4409
	/gene="DDX11L1"
	/note="DEAD/H (Asp-Glu-Ala-Asp/His) box helicase 11 like 1: Derived by automated computational analysis using gene prediction method: BestRefSeq."
	/pseudo
	/db_xref="GeneID:100287102"
	/db_xref="HGNC:37102"
misc_BNA	join(1874..2227,2613..2721,3221..4409)
	/gene="DDX11L1"
	/product="DEAD/H (Asp-Glu-Ala-Asp/His) box helicase 11 like 1"
	/note="Derived by automated computational analysis using gene prediction method: BestRefSeq."
	/pseudo
	/transcript_id="NR_046018.2"
	/db_xref="GI:389886562"
	/db_xref="GeneID:100287102"
	/db_xref="HGNC:37102"
	complement(4362..19370)
gene	/gene="WASH7P"
	/gene_synonym="FAM39F; WASH5P"
	/note="WAS protein family homolog 7 pseudogene: Derived by automated computational analysis using gene prediction method: BestRefSeq."
	/pseudo
	/db_xref="GeneID:653635"
	/db_xref="HGNC:38034"
	complement(join(<4362..4829,4970..5038,5796..5947,6507..6785,6858..7055,7233..7368,7606..7742,7915..8061,8268..8366,14738..14893,15232..15370))
	/gene="WASH7P"
	/gene_synonym="FAM39F; WASH5P"
	/product="WAS protein family homolog 7 pseudogene"
	/note="Derived by automated computational analysis using gene prediction method: BestRefSeq."
	/pseudo
	/transcript_id="NR_024540.1"
	/db_xref="GI:215277009"
	/db_xref="GeneID:653635"

```

STS      4596..4719      Sequence tagged site
/standard_name="ST13109350"
/db_xref="dbSTS:109350"
20366..20503
gene     /gene="MIR1302-2"
/gene_synonym="hsa-mir-1302-2; MIR1302-2"
/notes="microRNA 1302-2; Derived by automated computational
analysis using gene prediction method: BestRefSeq."
/db_xref="GeneID:100302278"
/db_xref="HGNC:35294"
precursor_RNA
20366..20503
/gene="MIR1302-2"
/product="microRNA 1302-2"
/notes="Derived by automated computational analysis using
gene prediction method: BestRefSeq."
/transcript_id="NM_01005484.1"
/db_xref="GeneID:100302278"
/db_xref="HGNC:35294"
/db_xref="MI0006363"
20438..20458
/gene="MIR1302-2"
/gene_synonym="hsa-mir-1302-2; MIR1302-2"
(ncRNA_class="miRNA"
/product="hsa-mir-1302"
/notes="Derived by automated computational analysis using
gene prediction method: BestRefSeq."
/db_xref="MI0006363"
/db_xref="GeneID:100302278"
/db_xref="HGNC:35294"
/db_xref="MI0006363"

```

```

gene     59091..60008      A protein-coding gene
/gene="OR4F5"
/notes="Olfactory receptor, family 4, subfamily F, member
5; Derived by automated computational analysis using gene
prediction method: BestRefSeq."
/db_xref="GeneID:79501"
/db_xref="HGNC:14825"
/db_xref="HPRD:14974"
59091..60008
/gene="OR4F5"
/product="Olfactory receptor, family 4, subfamily F,
member 5"
/notes="Derived by automated computational analysis using
gene prediction method: BestRefSeq."
/transcript_id="NM_001005484.1"
/db_xref="GeneID:79501"
/db_xref="HGNC:14825"
/db_xref="HPRD:14974"
59091..60008
/gene="OR4F5"
/notes="Derived by automated computational analysis using
gene prediction method: BestRefSeq."
/codon_start=1
/product="Olfactory receptor 4F5"
/protein_id="NP_001005484.1"
/db_xref="GI:53828740"
/db_xref="CCDS:CCDS30547.1"
/db_xref="GeneID:79501"
/db_xref="HGNC:14825"
/db_xref="HPRD:14974"
/translation="MVTETIFLGLSDSGLQTLPLMFPVYGGIVGNLLIVTVVS
DRLGLDHPYVPLALMLIDLSSVYPMIMTDYFQGRVRSYKGLVQVPLDLPFGQ
SRWVLLAWGFRVIACTPRTTLMGRCVGLMAYWVWVQVQVLAFAVQLL
FCQSNVDSYVCDLPRVILACTDLYRLDMVANSGLVCSFVLLIISYTIIMTI
QHRPLKRSKALSTLNTVTVLLFPQVYVYANPPFISLDRFLAVYSVITPLAN
P1YVLRQWMTALRQLRQWNSVYV"

```

```

gene     complement(132447..164392)
/gene="LOC100996442"
/notes="Derived by automated computational analysis using
gene prediction method: Genom. "
/db_xref="GeneID:100996442"
A gene with 2 splice variants
misc_RNA complement(join(132447..133011,145767..145831,
154263..154791,155884..155942,158100..158165,
163753..164392))
/gene="LOC100996442"
/product="uncharacterized LOC100996442, transcript variant
x2"
/notes="Derived by automated computational analysis using
gene prediction method: Genom. Supporting evidence
includes similarity to: 1 EST, and 88% coverage of the
annotated genomic feature by RNAseq alignments, including
3 samples with support for all annotated introns"
/transcript_id="XR_246629.1"
/db_xref="GI:530360230"
misc_RNA complement(join(148390..148674,154263..154791,
155884..155942,158100..158165,163753..164392))
/gene="LOC100996442"
/product="uncharacterized LOC100996442, transcript variant
x1"
/notes="Derived by automated computational analysis using
gene prediction method: Genom. Supporting evidence
includes similarity to: 4 ESTs, and 86% coverage of the
annotated genomic feature by RNAseq alignments, including
4 samples with support for all annotated introns"
/transcript_id="XR_159064.3"
/db_xref="GI:530360229"
/db_xref="GeneID:100996442"

```

```

assembly_gap 167418..217417      Gap in the sequence
/estimated_length=50000
/gap_type="within scaffold"
/linkage_evidence="unspecified"
complement(217770..218778)
/gene="RPL23A2P1"
Another pseudogene
/gene_synonym="RPL23a_1.1"
/notes="ribosomal protein L23a pseudogene 21; Derived by
automated computational analysis using gene prediction
method: Curated Genom. "
/gene=do
/db_xref="GeneID:728481"
/db_xref="HGNC:35827"
COR11G join(complement(AP006221..136117..136731),AL627309.15:103..166904,
gap(50000),AP006222.1:1..40302)
//
• Note that instead of the nucleotide sequences
there is instead a reference to several
sequences joined together and a gap
• Reflects that the "completed genome" is anything
but!
```

NCBI genomes: Bioprojects

- Allow multiple related data to be linked together

Organism Overview: [Genome Project Report](#) | [Genome Annotation Report](#) | [Plasmid Annotation Report](#)

Streptomyces hygroscopicus

Streptomyces hygroscopicus overview

Lineage: Bacteria[6088] Actinobacteria[723] Actinobacteriales[4733] Actinobacteriales[977] Actinomycetales[347] Streptomycetales[36] Streptomycetaceae[302] Streptomyces[403] Streptomyces hygroscopicus[41]

Streptomyces. These bacteria are widely distributed in nature, especially in the soil. The characteristic earthy smell of freshly plowed soil is actually attributed to the aromatic bicyclic sesquiterpene produced by species of Streptomyces. There are currently 364 known species of this genus, many of which are the most important industrial *fungi*...

Representative genome: [View all sequenced](#)
 Streptomyces hygroscopicus subsp. *hygroscopicus* 5008

Genome Sequencing Projects

Organism	Bioproject	Assembly	Status	Chr	Plasmids	Size (Mb)	GC%	Gene	Protein
Streptomyces hygroscopicus subsp. <i>hygroscopicus</i> 5008	PRJNA8495	ALUJG0155.1	●	1	2	10.58	71.8	8,184	8,187
Streptomyces hygroscopicus subsp. <i>hygroscopicus</i> TUS2	PRJNA1087	ABG00849.1	●	1	2	10.58	71.8	8,064	8,017

NCBI genomes: Assembly

ASM24555v1

Organism name: *Streptomyces hygroscopicus* subsp. *hygroscopicus* 5008

Submitter: State Key Laboratory of Microbial Metabolism, Shanghai Jiao Tong University

Date: 20150107

Assembly level: *De novo* Chromosome

Genome representation: full

Reference Assembly ID: GCF_0024555.1 (draft)

Reference Assembly and GenBank Assembly Identical: yes

History: [View History](#)

Global statistics

Total sequence length	10,583,954
Total assembly gap length	0
Total number of chromosomes and plasmids	3

Assembly Definition: [View Statistics](#)

Global assembly definition

Assembly Unit: Primary Assembly (GCF_0024555.1)

Unit name	GenBank ID	RefSeq ID
Chromosome	GSS01372.1	NC_011796.1
Plasmid pJW101	GSS01373.1	NC_011796.1
Plasmid pJW102	GSS01374.1	NC_011796.1

NCBI genomes

- Frankly, this is a bit of a mess on the website
 - Often, records are not in their complete form, i.e., instead consisting of links to subordinate data types, e.g., assembly->nucleotide->protein
 - But sometimes these are from GenBank, other times from Refseq
- Genbank, Refseq, Entrez data structures to not overlap in a simple manner
- If you want to download genomic data, this is not the way to do it!

NCBI FTP portal

- <ftp://ftp.ncbi.nlm.nih.gov/>
- Allows access both from the web and from Terminal (using the program `ftp`)
- Includes all of NCBI's data structures as [separate](#) database folders

Index of <ftp://ftp.ncbi.nlm.nih.gov/>

Up to higher level directory

Name	Size	Last Modified
1000genomes		2/26/2013 12:00:00 AM
1000G	102405024 KB	11/15/2013 6:44:00 PM
1000G2	102405024 KB	8/30/2013 6:19:00 PM
1000G2.2	102405024 KB	8/30/2013 6:19:00 PM
1000G2.3	2 KB	8/28/2013 9:37:00 PM
1000G2.4		8/28/2013 9:30:00 PM
1000G2.5	12311/2013 8:55:00 PM	
1000G2.6	1/26/2014 4:10:00 PM	
1000G2.7	1/26/2014 11:55:00 AM	
1000G2.8	5/24/2012 12:00:00 AM	
1000G2.9	9/13/2004 12:00:00 AM	
1000G2.10	12/5/2013 4:31:00 AM	
1000G2.11	11/6/2013 7:00:00 PM	
1000G2.12	7/18/2013 12:00:00 AM	
1000G2.13	10/13/2011 12:00:00 AM	
1000G2.14	9/26/2013 7:39:00 PM	
1000G2.15	8/4/2006 1:00:00 AM	
1000G2.16	8/28/2013 9:38:00 PM	
1000G2.17	12/15/2013 9:13:00 PM	
1000G2.18	1/14/2014 8:07:00 PM	
1000G2.19	1/28/2014 2:01:00 PM	
1000G2.20	1/28/2014 2:01:00 PM	
1000G2.21	9/20/2011 12:00:00 AM	
1000G2.22	5/15/2013 12:00:00 AM	
1000G2.23	12/13/2013 9:18:00 PM	

Index of <ftp://ftp.ncbi.nlm.nih.gov/genomes/>

Up to higher level directory

Name	Size	Last Modified
AC159810.1_BAC1250		1/21/2014 4:30:00 PM
AC159810.1_BPC1251		1/20/2014 9:32:00 AM
Acyrthosiphon_pisum		5/15/2012 12:00:00 AM
Ades_argyph		9/21/2010 12:00:00 AM
Akaraopoda_melanosticta		8/6/2012 12:00:00 AM
Albugo_taroensis		12/10/2013 2:31:00 PM
Albugo_taroensis		12/7/2013 2:22:00 PM
Alphabacterium_graminis		5/15/2012 12:00:00 AM
Ames_glycophytocola		10/29/2013 5:52:00 PM
Amoeba_cuculianus		3/15/2012 12:00:00 AM
Amphiprion_nano		1/4/2008 12:00:00 AM
Amphiprion_nano		1/6/2014 12:00:00 AM
Amphiprion_nano		4/15/2012 12:00:00 AM
Amphiprion_nano		1/6/2014 4:42:00 PM
Amphiprion_nano		7/24/2013 12:00:00 AM
Amphiprion_nano		5/15/2012 12:00:00 AM
Amphiprion_nano		9/20/2008 12:00:00 AM
Amphiprion_nano		5/15/2010 12:00:00 AM
Amphiprion_nano		1/28/2014 1:29:00 PM
Amphiprion_nano		4/29/2013 12:00:00 AM
Amphiprion_nano		5/15/2012 12:00:00 AM
Amphiprion_nano		5/15/2012 12:00:00 AM
Amphiprion_nano		7/8/2013 12:00:00 AM
Amphiprion_nano		11/4/2013 7:02:00 PM
Amphiprion_nano		8/8/2013 7:01:00 PM

Index of <ftp://ftp.ncbi.nlm.nih.gov/genomes/Bacteria/>

Up to higher level directory

Name	Size	Last Modified
AF1900000000.1.11117.tig	1 KB	2/27/2013 12:00:00 AM
AAH000000000.1.11222.tig	1 KB	2/27/2013 12:00:00 AM
Alcalyculon_mariae_MBC11017_uid51617		6/12/2013 12:00:00 AM
Aerobacter_pasteurianus_3088_uid514193		8/5/2013 4:07:00 PM
Aerobacter_pasteurianus_PO_3383_01_47C_uid518377		6/12/2013 12:00:00 AM
Aerobacter_pasteurianus_PO_3383_01_uid499279		6/12/2013 12:00:00 AM
Aerobacter_pasteurianus_PO_3383_01_uid518373		6/15/2013 12:00:00 AM
Aerobacter_pasteurianus_PO_3383_01_uid518381		6/15/2013 12:00:00 AM
Aerobacter_pasteurianus_PO_3383_01_uid518379		6/15/2013 12:00:00 AM
Aerobacter_pasteurianus_PO_3383_01_uid518383		6/15/2013 12:00:00 AM
Aerobacter_pasteurianus_PO_3383_01_uid518381		6/15/2013 12:00:00 AM
Aerobacter_pasteurianus_PO_3383_01_uid518375		6/15/2013 12:00:00 AM
Aerobacterium_wobeli_DSM_14330_uid88073		6/15/2013 12:00:00 AM
Aerobacterium_undatum_DSM_5501_uid81423		6/15/2013 12:00:00 AM
Aerobacterium_undatum_PO_34_uid80812		6/15/2013 12:00:00 AM
Aerobacterium_undatum_AJ_uid80809		6/15/2013 12:00:00 AM
Aerobacterium_undatum_NBCC_15126_uid232761		11/20/2013 5:12:00 AM
Aerobacterium_undatum_uid80813		8/7/2013 4:13:00 AM
Aeromonas_coccoloba_fermentans_DSM_20731_uid81471		6/15/2013 12:00:00 AM
Aeromonas_coccoloba_fermentans_PO_34_uid80812		6/15/2013 12:00:00 AM
Aeromonas_hydrophila_UID_uid86675		6/15/2013 12:00:00 AM

Index of ftp://ftp.ncbi.nlm.nih.gov/genomes/Bacteria/Acholeplasma_laidlawii_PG_8A_uid58901/

Up to higher level directory

Name	Size	Last Modified
NC_010163.GenBank-2.5m	343 KB	1/6/2011 12:00:00 AM
NC_010163.GenBankMM-2.6	83 KB	1/6/2011 12:00:00 AM
NC_010163.GenBank	15 KB	1/6/2011 12:00:00 AM
NC_010163.Prodigal-2.50	296 KB	1/6/2011 12:00:00 AM
NC_010163.kan	3660 KB	6/12/2013 12:00:00 AM
NC_010163.faa	579 KB	12/29/2012 12:00:00 AM
NC_010163.gff	1485 KB	8/27/2012 12:00:00 AM
NC_010163.fna	1482 KB	1/6/2011 12:00:00 AM
NC_010163.frm	15 KB	10/12/2010 12:00:00 AM
NC_010163.gbk	4783 KB	6/12/2013 12:00:00 AM
NC_010163.gzi	3 KB	6/12/2013 12:00:00 AM
NC_010163.gff	138 KB	12/29/2012 12:00:00 AM
NC_010163.gff	111 KB	1/6/2011 12:00:00 AM
NC_010163.msi	3 KB	1/6/2011 12:00:00 AM
NC_010163.gzi	1 KB	12/29/2012 12:00:00 AM
NC_010163.fna	1 KB	6/12/2013 12:00:00 AM
NC_010163.fna	1354 KB	6/12/2013 12:00:00 AM

.fna: chromosome nucleotides
 .ffn: genes
 .faa: proteins
 .gbk: annotated chromosomes
 .gbs: scaffolds
 .gff: all annotations
 .pft: protein annotations
 .rnt: RNA annotations (not mRNA)

