

Computer lab exercises #8

Comments on projects worth sharing:

1. Use BLINK whenever possible. It can save a lot of waiting and greatly accelerates explorations.

From a protein sequence entry in NCBI select “BLINK” under related information. (You might need to scroll down, in case the upper tables are expanded).

BLINK provides a GUI interface to pre-computed BLASTP searches.

Display Settings: GenPept

Send to:

tumor necrosis factor alpha [Homo sapiens]

GenBank: ACO37640.1

[FASTA](#) [Graphics](#)

[Go to:](#)

LOCUS ACO37640 232 aa linear PRI 06-APR-2009
 DEFINITION tumor necrosis factor alpha, partial [Homo sapiens].
 ACCESSION ACO37640
 VERSION ACO37640.1 GI:226201421
 DBSOURCE accession [FJ795028.1](#)
 KEYWORDS .
 SOURCE Homo sapiens (human)
 ORGANISM [Homo sapiens](#)
 Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
 Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini;
 Catarrhini; Hominidae; Homo.
 REFERENCE 1 (residues 1 to 232)
 AUTHORS Guan,W.J., Ma,Y.H., Yu,L.L., Na,R.S. and Liu,S.
 TITLE Direct Submission
 JOURNAL Submitted (28-FEB-2009) Academy of Agricultural Sciences, Institute
 of Animal Science, Quanmingyuan West, Beijing 100193, People's
 Republic of China

FEATURES Location/Qualifiers
 source 1..232
 /organism="Homo sapiens"
 /db_xref="taxon:9606"
 /chromosome="6"
 /map="6p21.3"
 /sex="male"
 /tissue_type="placenta"
 /country="China"
 /collection_date="2008"
 <1..232
 /product="tumor necrosis factor alpha"
 /name="APC1 protein"
 87..230
 /region_name="TNF"
 /note="Tumor Necrosis Factor; TNF superfamily members
 include the cytokines: TNF (TNF-alpha), LT
 (lymphotoxin-alpha, TNF-beta), CD40 ligand, Apo2L (TRAIL),
 Fas ligand. and osteoprotegerin (OPG) ligand. These

[Protein](#)

[Region](#)

Change region shown

Customize view

Analyze this sequence

Protein 3D Structure

Articles about the TNF gene

Pathways for the TNF gene

Reference sequence information

More about the TNF gene

Homologs of the TNF gene

LinkOut to external resources

Related information

[BLink](#)

[Related Sequences](#)

[BioSystems](#)

[CDD Search Results](#)

[Conserved Domains \(Concise\)](#)

[Conserved Domains \(Full\)](#)

[Domain Relatives](#)



Pre-computed BLAST results for: [gi|226201421|gb|ACO37640.1](#) tumor necrosis factor alpha [Homo sapiens]

Total (score > 100) : 1804 hits in 1804 proteins in 205 species

Selected: 1804 hits in 1804 proteins in 205 species Filter: Min Score: 100 |

Other views (Reports): [Taxonomy report](#) [Multiple Alignment](#) [Blast](#)

[Reset all filters](#)

► [Choose Display Options](#)

Archaea Bacteria Metazoa Fungi Plants Viruses The Others [reset selection](#)

Results: 1 - 100 [Next Page](#) [Last](#)

SCORE	ACCESSION	Length	Protein Description
reset selection			
232 aa			
Conserved Domain Database hits			
1188	BAC54944	233	tumor necrosis factor [Homo sapiens]
1188	AAO21132	233	tumor necrosis factor (TNF superfamily, member 2) [Homo sapiens]
1188	BAG37464	233	unnamed protein product [Homo sapiens]
1188	EAX03424	233	tumor necrosis factor (TNF superfamily, member 2) [Homo sapiens]
1188	AAX41550	233	tumor necrosis factor [synthetic construct]
1188	CAA26669	233	TNF-alpha [Homo sapiens]
1188	CAA78745	233	tumor necrosis factor, Tnfa [Homo sapiens]
1188	CAA25650	233	unnamed protein product [Homo sapiens]
1188	P01375	233	RecName: Full=Tumor necrosis factor; AltName: Full=Cachectin; AltName: Full=TNF-a
1188	gi 224323	233	tumor necrosis factor
1188	gi 224436	233	tumor necrosis factor
1188	AAA61200	233	tumor necrosis factor [Homo sapiens]
1188	AAA36758	233	tumor necrosis factor precursor [Homo sapiens]
1188	BAF31279	233	TNFA protein [Homo sapiens]
1188	ABM82588	233	tumor necrosis factor (TNF superfamily, member 2) [synthetic construct]
1188	ABM85775	233	tumor necrosis factor (TNF superfamily, member 2) [synthetic construct]
1188	XP_003831637	233	PREDICTED: tumor necrosis factor [Pan paniscus]
1188	BAG73840	233	tumor necrosis factor [synthetic construct]
1188	AHJ25918	233	tumor necrosis factor [Homo sapiens]
1188	AHJ25919	233	tumor necrosis factor [Homo sapiens]
1188	AHJ25920	233	tumor necrosis factor [Homo sapiens]
1188	AHJ25921	233	tumor necrosis factor [Homo sapiens]
1188	AHJ25922	233	tumor necrosis factor [Homo sapiens]
1188	AHJ25923	233	tumor necrosis factor [Homo sapiens]

▼ Choose Display Options



filter hits

best hits all hits hide identical ⓘ

Minimum Hit Score

100

Maximum Hit Score

New Search By GI

GO

Items per page

100

BLINK

Parameters have been changed. Please, press BLINK button to update the view.

Pre-computed BLAST results for: [gi|226201421|gb|ACO37640.1](#) tumor necrosis factor alpha [Homo sapiens]

Total (score > 100) : 1804 hits in 1804 proteins in 205 species

Selected: 1804 hits in 1804 proteins in 205 species Filter: Min Score: 100 |

Other views (Reports):

[Reset all filters](#)

► [Choose Display Options](#)

Archaea Bacteria Metazoa Fungi Plants Viruses The Others [reset selection](#)

Results: 1 - 100 [Next Page](#) [Last](#)

blink	SCORE	ACCESSION	Length	Protein Description
				Conserved Domain Database hits
◆	1188	BAC54944	233	tumor necrosis factor [Homo sapiens]
◆	1188	AAO21132	233	tumor necrosis factor (TNF superfamily, member 2) [Homo sapiens]
◆	1188	BAG37464	233	unnamed protein product [Homo sapiens]
◆	1188	EAX03424	233	tumor necrosis factor (TNF superfamily, member 2) [Homo sapiens]
◆	1188	AAX41550	233	tumor necrosis factor [synthetic construct]
◆	1188	CAA26669	233	TNF-alpha [Homo sapiens]
◆	1188	CAA78745	233	tumor necrosis factor, Tnfa [Homo sapiens]
◆	1188	CAA25650	233	unnamed protein product [Homo sapiens]
◆	1188	P01375	233	RecName: Full=Tumor necrosis factor; AltName: Full=Cachectin; AltName: Full=TNF-a
◆	1188	gi 224323	233	tumor necrosis factor
◆	1188	gi 224436	233	tumor necrosis factor
◆	1188	AAA61200	233	tumor necrosis factor [Homo sapiens]
◆	1188	AAA36758	233	tumor necrosis factor precursor [Homo sapiens]
◆	1188	BAF31279	233	TNFA protein [Homo sapiens]
◆	1188	ABM82588	233	tumor necrosis factor (TNF superfamily, member 2) [synthetic construct]
◆	1188	ABM85775	233	tumor necrosis factor (TNF superfamily, member 2) [synthetic construct]
◆	1188	XP_003831637	233	PREDICTED: tumor necrosis factor [Pan paniscus]
◆	1188	BAG73840	233	tumor necrosis factor [synthetic construct]
◆	1188	AHJ25918	233	tumor necrosis factor [Homo sapiens]
◆	1188	AHJ25919	233	tumor necrosis factor [Homo sapiens]
◆	1188	AHJ25920	233	tumor necrosis factor [Homo sapiens]
◆	1188	AHJ25921	233	tumor necrosis factor [Homo sapiens]
◆	1188	AHJ25922	233	tumor necrosis factor [Homo sapiens]
◆	1188	AHJ25923	233	tumor necrosis factor [Homo sapiens]

Pre-computed BLAST results for: [gi|25952111|ref|NP_000585.2](#) tumor necrosis factor [Homo sapiens]

Matching gis: [27544420;27802685;510158303;363836827;189054614;118767860;119623829;375725640;383266295;380291495](#)

Total (score > 100) : 1017 hits in 1017 proteins in 186 species

Selected: 1017 hits in 1017 proteins in 186 species Filter: Min Score: 100 | Best hits |

Other views (Reports): [Taxonomy report](#) [Multiple Alignment](#) [Blast](#)

[Reset all filters](#)

► [Choose Display Options](#)

Archaea Bacteria 1002 Metazoa Fungi Plants Viruses 15 The Others [reset selection](#)

Results: [First](#) [Previous Page](#) 101 - 186

% hits



233 aa

[reset selection](#)

blink



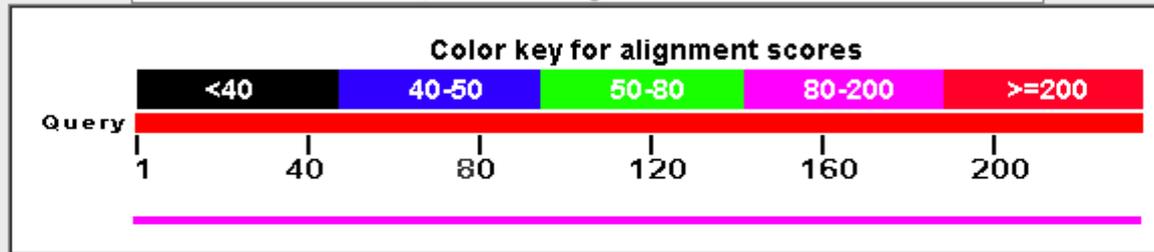
SCORE ACCESSION N Tax

[Conserved Domain Database hits](#)

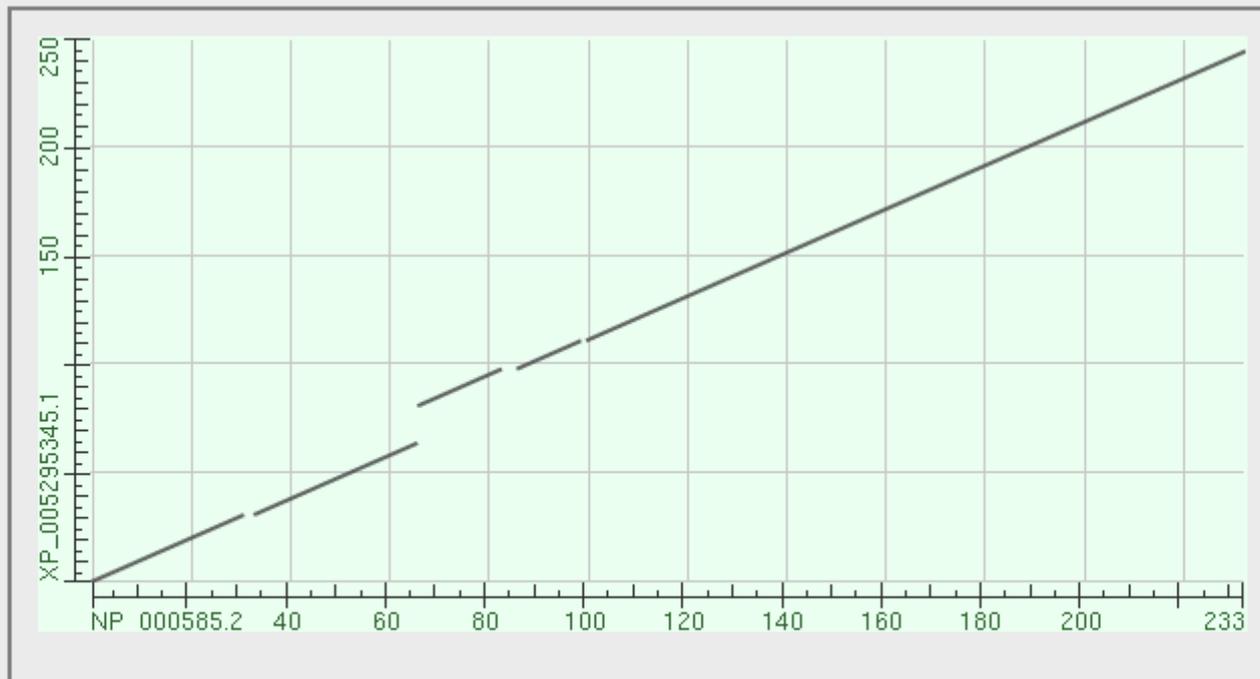
◆	—	484	CAD80055	1	Elephas maximus
◆	—	478	CAD80061	1	Trichys fasciculata
◆	—	466	CAD80060	1	Castor canadensis
◆	—	462	CAD80059	1	Anomalurus sp. T-1787
◆	—	460	CAD80057	1	Dipus sagitta
◆	—	450	XP_005295345	10	Chrysemys picta bellii
◆	—	423	CAD80058	1	Dipodomys merriami
◆	—	392	XP_003218070	7	Anolis carolinensis
◆	—	390	ABC88246	1	Phodopus sungorus
◆	—	285	NP_001108250	3	Xenopus laevis
◆	—	274	BAB68749	1	Mus musculus brevisrostris
◆	—	266	NP_001107143	11	Xenopus (Silurana) tropicalis
◆	—	258	AAP94278	1	Acanthopagrus schlegelii
◆	—	257	AGQ17907	1	Ginglymostoma cirratum

Distribution of 1 Blast Hits on the Query Sequence

Mouse over to see the define, click to show alignments



Plot of [gj|25952111|ref|NP_000585.2](#) vs [gj|530607321|ref|XP_005295345.1](#)



If one is looking for homologs in other phyla or domains,

Pre-computed BLAST results for: [gi|16764351|ref|NP_459966.1](#) metallothionein SmtA [Salmonella enterica subsp. enterica serovar Typhimurium str. LT2]

Matching gis: [353075391;374353981;194406803;194448499;194459024;514603770;514609876;514610818;514618890;514627615;514631814;514637706;514641511;514641512](#)

Total (score > 100) : 26421 hits in 26413 proteins in 8653 species

Selected: 26421 hits in 26413 proteins in 8653 species Filter: Min Score: 100 |

Other views (Reports): [Taxonomy report](#) [Multiple Alignment](#) [Blast](#)

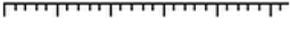
[Reset all filters](#)

[Choose Display Options](#)

[547](#) [Archaea](#) [25115](#) [Bacteria](#) [40](#) [Metazoa](#) [217](#) [Fungi](#) [30](#) [Plants](#) [2](#) [Viruses](#) [470](#) [The Others](#) [reset selection](#)

Results: 1 - 100 [Next Page](#) [Last](#)

% hits  [reset selection](#)
267 aa

blink	SCORE	ACCESSION	Length	Protein Description
				Conserved Domain Database hits
	1405	EHB41151	267	methyltransferase domain protein [Salmonella enterica subsp. enterica serovar
	1405	AEZ45742	267	hypothetical protein STBHUCCB_20590 [Salmonella enterica subsp. enterica serov
	1405	ACF67022	267	SmtA protein [Salmonella enterica subsp. enterica serovar Heidelberg str. SL47
	1405	YP_002044983	267	metallothionein SmtA [Salmonella enterica subsp. enterica serovar Heidelberg s
	1405	EDX47863	267	SmtA protein [Salmonella enterica subsp. enterica serovar Kentucky str. CVM291
	1405	EPI63832	267	methyltransferase domain protein [Salmonella enterica subsp. enterica serovar
	1405	EPI69660	267	methyltransferase domain protein [Salmonella enterica subsp. enterica serovar
	1405	EPI70579	267	methyltransferase domain protein [Salmonella enterica subsp. enterica serovar
	1405	EPI78455	267	methyltransferase domain protein [Salmonella enterica subsp. enterica serovar
	1405	EPI86807	267	methyltransferase domain protein [Salmonella enterica subsp. enterica serovar
	1405	EPI90919	267	methyltransferase domain protein [Salmonella enterica subsp. enterica serovar
	1405	EPI96476	267	methyltransferase domain protein [Salmonella enterica subsp. enterica serovar

Pre-computed BLAST results for: [gi|16764351|ref|NP_459966.1](#) metallothionein SmtA [Salmonella enterica subsp. enterica serovar Typhimurium str. LT2]

Matching gis: [353075391;374353981;194406803;194448499;194459024;514603770;514609876;514610818;514618890;514627615;514631814;514637706;514641511;51](#)

Total (score > 100) : 26421 hits in 26413 proteins in 8653 species

Selected: 40 hits in 40 proteins in 27 species Filter: Min Score: 100 | Included taxons: Metazoa; |

Other views (Reports): [Taxonomy report](#) [Multiple Alignment](#) [Blast](#)

[Reset all filters](#)

▶ [Choose Display Options](#)

Archaea Bacteria **40** Metazoa Fungi Plants Viruses The Others [reset selection](#)

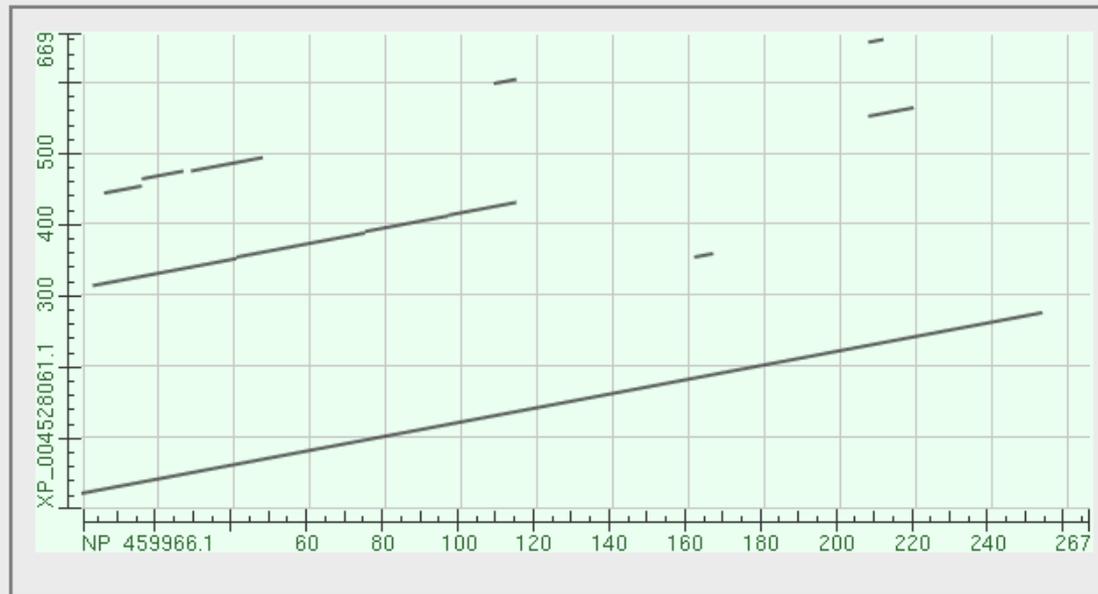
Results: 1 - 40

% hits [reset selection](#)
267 aa

blink

	SCORE	ACCESSION	Length	Protein Description	
				Conserved Domain Database hits	
◆	1109	XP_004528061	669	PREDICTED: chromosome partition protein MukF-like [Ceratitis capitata]	
◆	935	GAA57985	361	3-demethylubiquinone-9 3-methyltransferase, partial [Clonorchis sinensis]	
◆	4430	XP_004532210	249	PREDICTED: 3-demethylubiquinone-9 3-methyltransferase-like [Ceratitis capitata]	
◆	467	XP_001603088	289	PREDICTED: hexaprenyldihydroxybenzoate methyltransferase, mitochondrial-like [
◆	207	123	CAB16512	268	Protein COQ-3, isoform a [Caenorhabditis elegans]
◆	116	123	NP_001041045	268	Protein COQ-3, isoform a [Caenorhabditis elegans]

Plot of gi|16764351|ref|NP_459966.1| vs gi|498975409|ref|XP_004528061.1|



PREDICTED: chromosome partition protein MukF-like [Ceratitis capitata]

Sequence ID: [ref|XP_004528061.1|](#) Length: 669 Number of Matches: 7

Range 1: 23 to 276 [GenPept](#) [Graphics](#)

▼ Next Match ▲ Previous Match

	Score	Expect	Identities	Positives	Gaps
	431 bits(1109)	1e-151	209/254(82%)	225/254(88%)	0/254(0%)
Query	1	MQDRNFDDIAEKFSRNIYGTTKQQLRQAILWQDLDRVLEEIGGRKLRVLDAGGGEGQTAI	60		
Sbjct	23	+QDRNFDDIAEKFSRNIYGTTKQQLRQAILWQDLDR+L G LR+LDAGGG GQTAI	82		
Query	61	KMAERGHQVTLCDLSGEMIARARQAAEAKGVSKDMHFIQCPAQDVASHLESPVDLILFHA	120		
Sbjct	83	+MAERGH VTLCDLS EMIA A++AA+ KGVS MHF+QC QDVA HLESPVDLILFHA	142		
Query	121	VLEWVADPVGVLETLWSVLRPGGALS LMFYNANGLLMHNMVAGNFYVQAGMPKRRKRTL	180		
Sbjct	143	VLEWVAEPRTVLDTLWSTLRPGGALS LMFYNANGLLHNMVATNFYVQAGMPKRRKRTL	202		
Query	181	SPDYPRDPAQVYQWLEAIGWQITGKTGVRVFDHDLREKHQQRDCYETLVELETRYCRQEP	240		
Sbjct	203	SPDYPRDP QVY WL+ GWQITGKTGVRVFDHDLREK +QRD Y L+ELETRYCRQEP	262		
Query	241	YISLGRYIHVTAIK	254		
Sbjct	263	+ISLGRYIHVTA K	276		

Comments on projects worth sharing:

2. The Taxonomy Browser at NCBI can facilitate finding genome and EST projects.

In entrez select taxonomy as databank



- Recent
- Protein
- ✓ PubMed
- All
- All Databases
- Assembly
- BioProject
- BioSample
- BioSystems
- Books
- ClinVar
- Clone
- Conserved Domains
- dbGaP
- dbVar
- Epigenomics
- EST
- Gene
- Genome
- GEO DataSets
- GEO Profiles
- GSS
- HomoloGene
- MedGen
- MeSH
- NCBI Web Site
- NLM Catalog
- Nucleotide
- OMIA
- OMIM
- PMC
- PopSet
- Probe
- Protein
- Protein Clusters
- PubChem BioAssay
- PubChem Compound
- PubChem Substance
- PubMed
- PubMed Health
- SNP
- SRA
- Structure
- Taxonomy
- ToolKit
- ToolKitAll
- ToolKitBook
- UniGene
- UniSTS

Search **Search** [Help](#)



more than 23 million citations for biomedical
NE, life science journals, and online books.
links to full-text content from PubMed Central
ES.

PubMed Commons

PubMed's new commenting system

[More](#)

- ### Using PubMed
- [PubMed Quick Start Guide](#)
 - [Full Text Articles](#)
 - [PubMed FAQs](#)
 - [PubMed Tutorials](#)
 - [New and Noteworthy](#)

- ### PubMed Tools
- [PubMed Mobile](#)
 - [Single Citation Matcher](#)
 - [Batch Citation Matcher](#)
 - [Clinical Queries](#)
 - [Topic-Specific Queries](#)

- ### More Resources
- [MeSH Database](#)
 - [Journals in NCBI Databases](#)
 - [Clinical Trials](#)
 - [E-Utilities](#)
 - [LinkOut](#)

You are here: NCBI > Literature

- ### GETTING STARTED
- [NCBI Education](#)
 - [NCBI Help Manual](#)
 - [NCBI Handbook](#)
 - [Training & Tutorials](#)

POPULAR	FEATURED	NCBI INFORMATION
PubMed	Genetic Testing Registry	About NCBI
Bookshelf	PubMed Health	Research at NCBI
PubMed Central	GenBank	NCBI News
PubMed Health	Reference Sequences	NCBI FTP Site
BLAST	Gene Expression Omnibus	NCBI on Facebook
Nucleotide	Map Viewer	NCBI on Twitter
Genome	Human Genome	NCBI on YouTube
SNP	Mouse Genome	
Gene	Influenza Virus	
Protein	Primer-BLAST	
PubChem	Sequence Read Archive	



Search for an organisms you are interested in,

NCBI Resources How To

Taxonomy

[Save search](#) [Limits](#) [Advanced](#)

[Display Settings:](#) Summary

 [Sepia](#)

genus, cephalopods

[Nucleotide](#)

[Protein](#)

Select the group of organisms you are interested in

NCBI Taxonomy Browser

Entrez PubMed Nucleotide Protein Genome Structure PMC Taxonomy

Search for as lock

Display levels using filter:

Nucleotide Nucleotide EST Nucleotide GSS Protein Structure Genome Popset SNP

Domains GEO Datasets UniGene UniSTS PubMed Central Gene HomoloGene SRA Experiments

MapView LinkOut BLAST TRACE Probe Assembly Bio Project Bio Sample

Bio Systems Clone DB dbVar Epigenomics GEO Profiles PubChem BioAssay Protein Clusters Host

[Lineage](#) (full): [root](#); [cellular organisms](#); [Eukaryota](#); [Opisthokonta](#); [Metazoa](#); [Eumetazoa](#); [Bilateria](#); [Protostomia](#); [Lophotrochozoa](#); [Mollusca](#); [Cephalopoda](#); [Coleoidea](#); [Neocoleoidea](#); [Decapodiformes](#); [Sepiida](#); [Sepiina](#); [Sepiidae](#)

o [Sepia](#) *Click on organism name to get more information.*

- [Sepia aculeata](#)
- [Sepia andreana](#)
- [Sepia apama](#) (giant Australian cuttlefish)
- [Sepia aureomaculata](#)
- [Sepia bertheloti](#)
- [Sepia elegans](#)
- [Sepia elliptica](#)
- [Sepia esculenta](#) (golden cuttlefish)
- [Sepia filibrachia](#)
- [Sepia furcata](#)
- [Sepia gibba](#)
- [Sepia hierredda](#)

Check genomes and ESTs and whatever else you like, THEN click display

NCBI Taxonomy Browser

Entrez PubMed Nucleotide Protein Genome Structure PMC Taxonomy

Search for [] as [complete name] [lock] [Go] [Clear]

Display [3] levels using filter: [none]

Nucleotide Nucleotide EST Nucleotide GSS Protein Structure Genome Popset SNP

Domains GEO Datasets UniGene UniSTS PubMed Central Gene HomoloGene SRA Experiments

MapView LinkOut BLAST TRACE Probe Assembly Bio Project Bio Sample

Bio Systems Clone DB dbVar Epigenomics GEO Profiles PubChem BioAssay Protein Clusters Host

Lineage (full): [root](#); [cellular organisms](#); [Eukaryota](#); [Opisthokonta](#); [Metazoa](#); [Eumetazoa](#); [Bilateria](#); [Protostomia](#); [Lophotrochozoa](#); [Mollusca](#)

- [Cephalopoda](#) (cephalopods) [114,034](#) [36](#) *Click on organism name to get more information.*
 - [Coleoidea](#) [112,997](#) [35](#)
 - [Neocoleoidea](#) [112,997](#) [35](#)
 - [Decapodiformes](#) [112,994](#) [24](#)
 - [Octopodiformes](#) [3](#) [11](#)
 - [Nautiloidea](#) [1,037](#) [1](#)
 - [Nautilida](#) [1,037](#) [1](#)
 - [Nautilidae](#) [1,037](#) [1](#)
 - [environmental samples](#)
 - [Cephalopoda environmental sample](#)

Disclaimer: The NCBI taxonomy database is not an authoritative source for nomenclature or classification - please consult the relevant scientific literature for the most reliable information.

If you move up or down the taxonomic hierarchy, the checked items will be displayed for each Taxon

NCBI Taxonomy Browser

Entrez PubMed Nucleotide Protein Genome Structure PMC

Search for as lock

Display levels using filter:

Nucleotide Nucleotide EST Nucleotide GSS Protein Structure Genome Popset

Domains GEO Datasets UniGene UniSTS PubMed Central Gene HomoloGene

MapView LinkOut BLAST TRACE Probe Assembly Bio Project

Bio Systems Clone DB dbVar Epigenomics GEO Profiles PubChem BioAssay Protein Clusters

Lineage (full): [root](#); [cellular organisms](#); [Eukaryota](#); [Opisthokonta](#); [Metazoa](#); [Eumetazoa](#); [Bilateria](#); [Protostomia](#); [Cephalopoda](#); [Coleoidea](#); [Neocoleoidea](#); [Decapodiformes](#)

○ [Sepioida](#) [44,499](#) [2](#) *Click on organism name to get more information.*

○ [Idiosepiidae](#) [9,079](#)

○ [Idiosepius](#) [9,079](#)

▪ [Idiosepius biserialis](#)

▪ [Idiosepius macrocheir](#)

▪ [Idiosepius notoides](#)

▪ [Idiosepius paradoxus](#) [9,079](#)

▪ [Idiosepius picteti](#)

▪ [Idiosepius pygmaeus](#)

▪ [Idiosepius thailandicus](#)

○ [Sepioidae](#) (bobtail squids) [35,420](#) [2](#) OR

○ [Euprymna](#) [35,420](#) [1](#)

▪ [Euprymna berryi](#)

▪ [Euprymna hyllebergi](#)

▪ [Euprymna morsei](#)

▪ [Euprymna scolopes](#) [35,420](#) [1](#)

▪ [Euprymna stenodactyla](#)

▪ [Euprymna tasmanica](#)

▪ [Euprymna sp.](#)

○ [Heteroteuthis](#)

▪ [Heteroteuthis dispar](#)

▪ [Heteroteuthis hawaiiensis](#)

▪ [Heteroteuthis rvukyuensis](#)

Retrieves all ESTs from Genbank

You can select format or send to file to get a multiple sequence fasta file with all sequences

NCBI Resources How To

EST EST txid34531[Organism:exp]
Save search Limits Advanced

Display Settings: Summary, 20 per page, Sorted by Default order Send to:

Format	Items per page	Sort by
<input checked="" type="radio"/> Summary	<input type="radio"/> 5	<input checked="" type="radio"/> Default order
<input type="radio"/> EST	<input type="radio"/> 10	<input type="radio"/> Accession
<input type="radio"/> GenBank	<input checked="" type="radio"/> 20	<input type="radio"/> Date Modified
<input type="radio"/> GenBank (full)	<input type="radio"/> 50	<input type="radio"/> Date Released
<input type="radio"/> FASTA	<input type="radio"/> 100	<input type="radio"/> Organism Name
<input checked="" type="radio"/> FASTA (text)	<input type="radio"/> 200	<input type="radio"/> Taxonomy ID
<input type="radio"/> ASN.1		
<input type="radio"/> Revision History		
<input type="radio"/> Accession List		
<input type="radio"/> GI List		

1 of 1771 Next > Last >>

[aeg-p-24-0-UI 3-, mRNA](#)

[aeg-p-22-0-UI 3-, mRNA](#)

731 bp linear mRNA
Accession: DW286721.1 GI: 84452125
[EST](#) [GenBank](#) [FASTA](#)

[UI-S-GU1-aeg-p-20-0-UI.s1 UI-S-GU1 Euprymna scolopes cDNA clone UI-S-GU1-aeg-p-20-0-UI 3-, mRNA sequence](#)

3. 638 bp linear mRNA
Accession: DW286720.1 GI: 84452124
[EST](#) [GenBank](#) [FASTA](#)

You can select format or send to file to get a multiple sequence fasta file with all sequences

www.ncbi.nlm.nih.gov/nucest/?term=txid34531[Organism:exp]

EST

Save search Limits Advanced

Display Settings: Summary, 20 per page, Sorted by Default order

Send to: **Filter your results:**

Choose Destination

File Clipboard
 Collections

Download 35420 items.

Format

Summary
 EST
 GenBank
 GenBank (full)
 FASTA
 ASN.1
 XML
 INSDSeq XML
 TinySeq XML
 Accession List
 GI List

Found 35908 nucleotide sequences. Nucleotide (488) EST (35420)

Results: 1 to 20 of 35420

1. [UI-S-GU1-aeg-p-24-0-UI.s1 UI-S-GU1 Euprymna scolopes cDNA clone UI-S-GU1-24-0-UI.s1](#)
[sequence](#)
696 bp linear mRNA
Accession: DW286722.1 GI: 84452126
[EST](#) [GenBank](#) [FASTA](#)

2. [UI-S-GU1-aeg-p-22-0-UI.s1 UI-S-GU1 Euprymna scolopes cDNA clone UI-S-GU1-22-0-UI.s1](#)
[sequence](#)
731 bp linear mRNA
Accession: DW286721.1 GI: 84452125
[EST](#) [GenBank](#) [FASTA](#)

3. [UI-S-GU1-aeg-p-20-0-UI.s1 UI-S-GU1 Euprymna scolopes cDNA clone UI-S-GU1-aeg-p-20-0-UI 3-, mRNA](#)
[sequence](#)
638 bp linear mRNA
Accession: DW286720.1 GI: 84452124
[EST](#) [GenBank](#) [FASTA](#)

Search details
txid34531[Organism:exp]

Short in class exercise

- Write and apply a script to calculate cumulative strand bias for a genome (fna file).

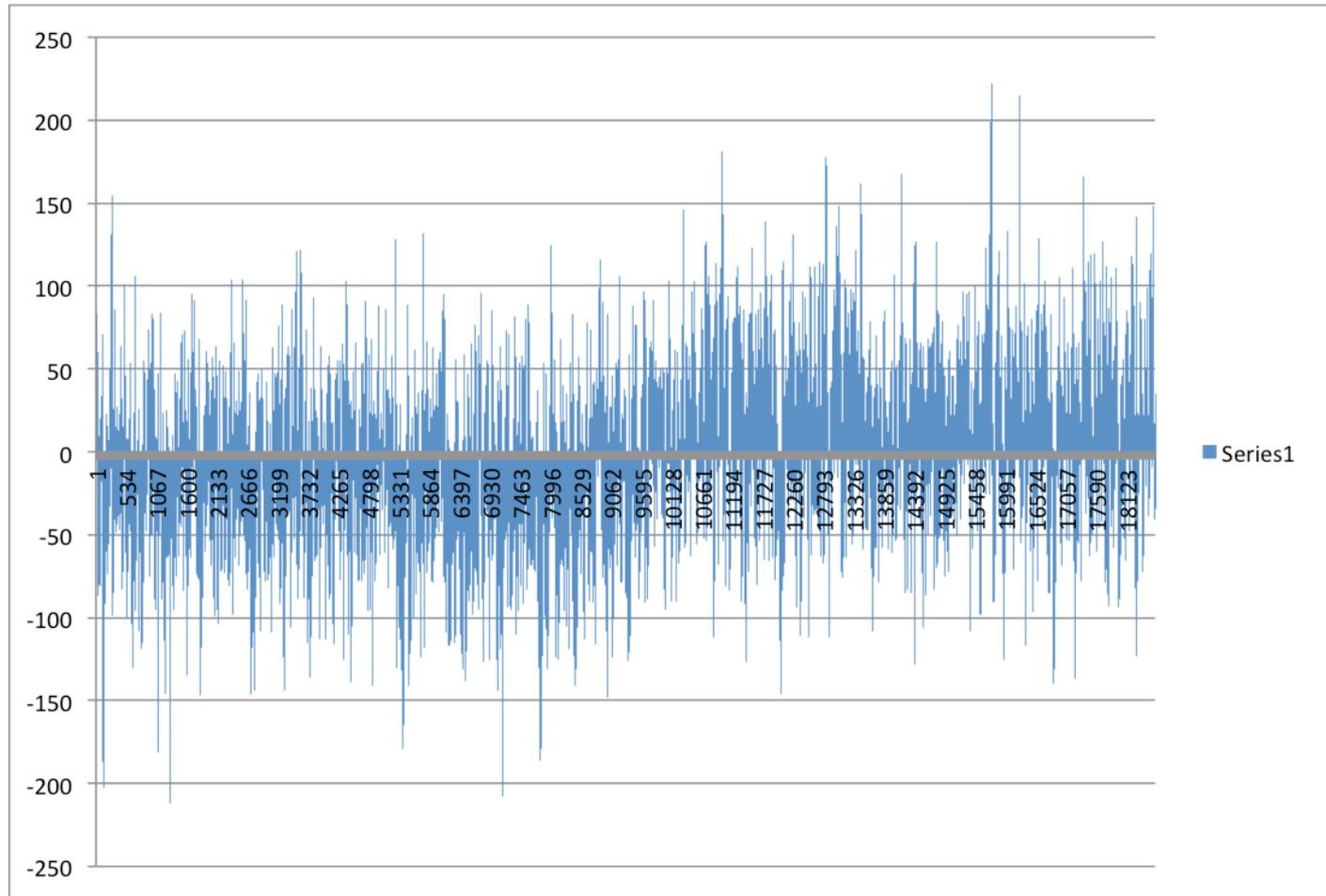
In DNA replication the two strands are not created equal.

See Drew Berry's Ted talk at <https://www.youtube.com/watch?v=WFCvkkDSfIU> for illustration (start at 3 min – 4.50 min, if not much time)

The differences between leading and lagging strand are reflected in the number of ORFs encoded on the strands and the presence of motifs that bind factors initiating and halting replication, and the composition with respect to nucleotides, and n-mers of nucleotides.

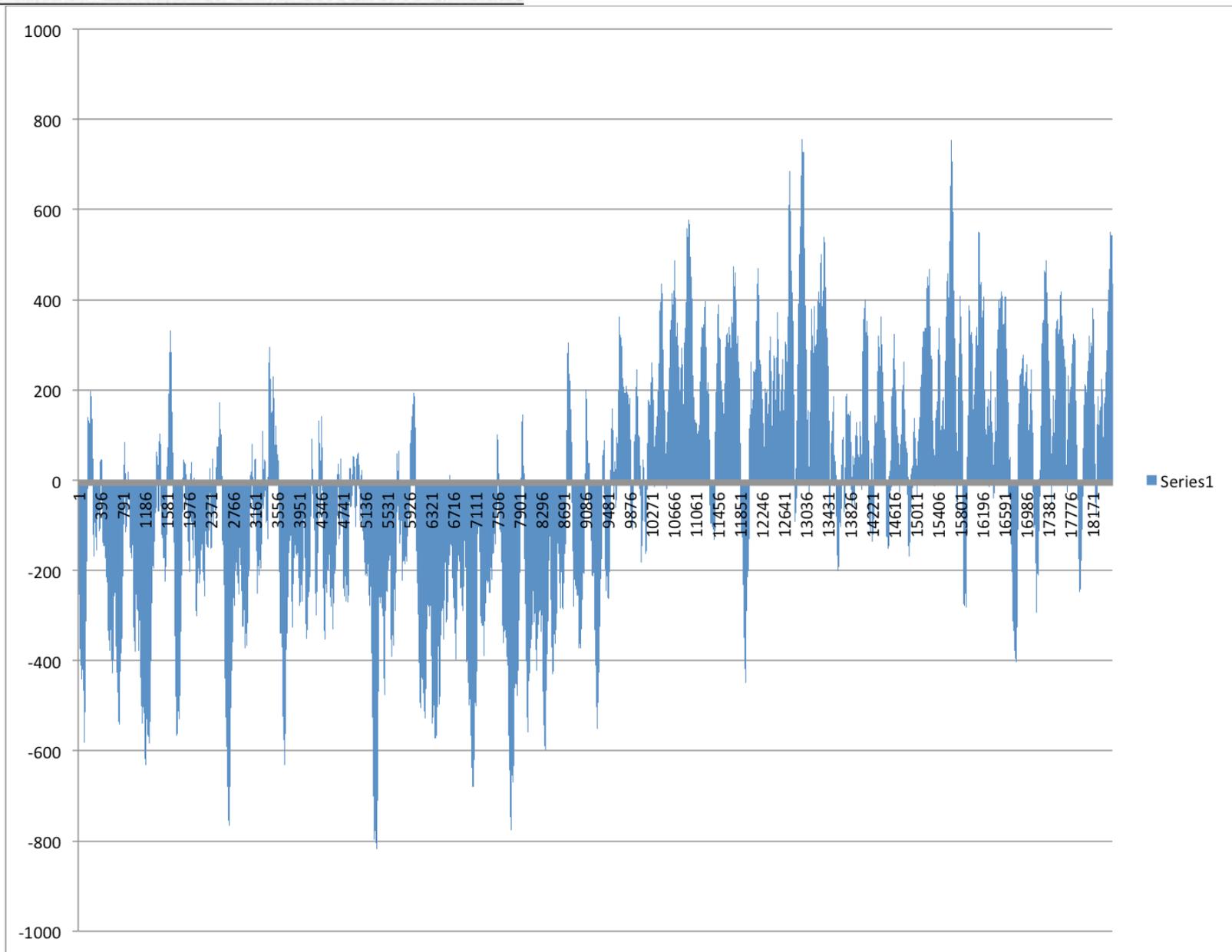
One way to look at strand bias, is to calculate the GC content in a rolling window.

Thermus thermophilus SG0.5JP17-16



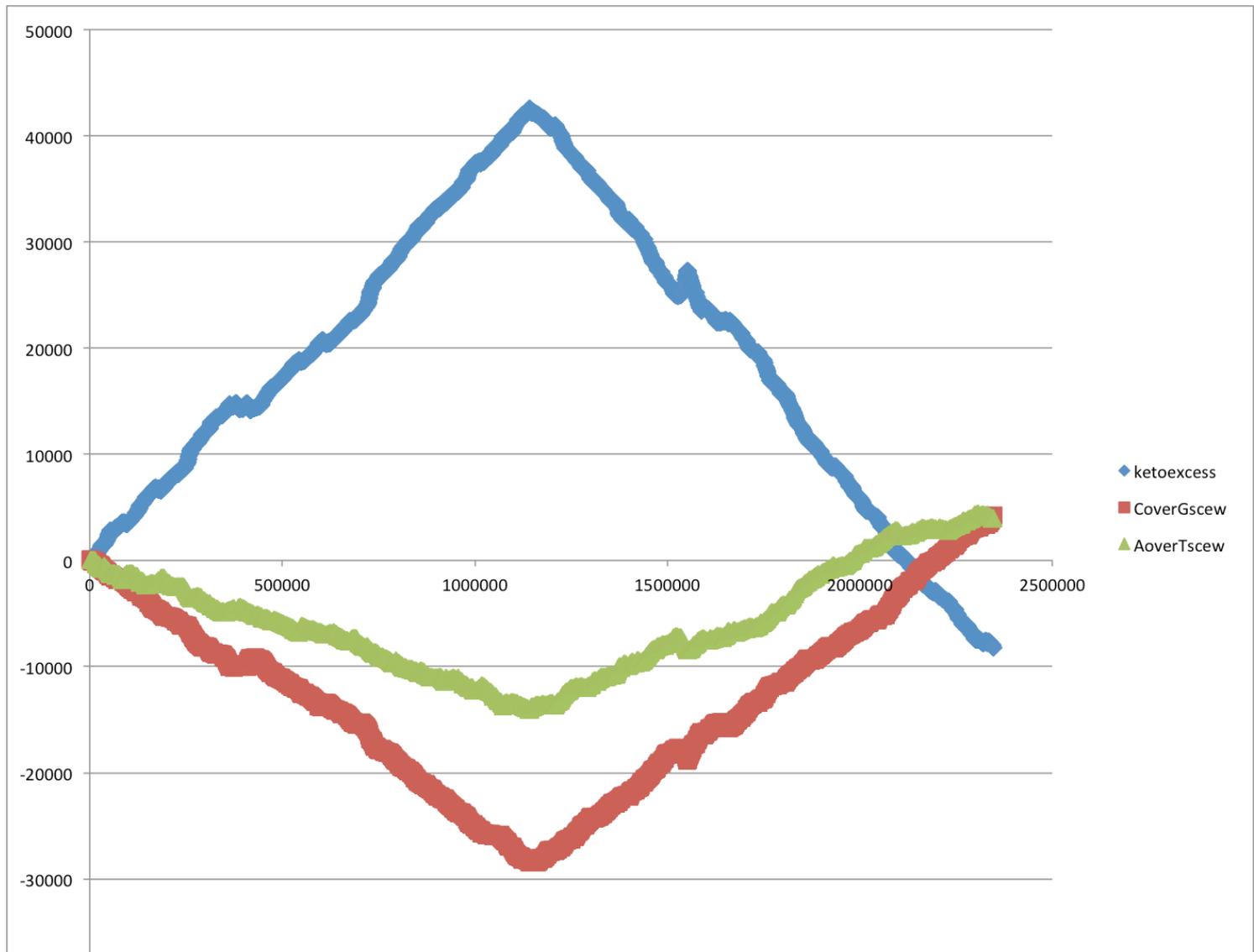
Window=1000 , printed every 100

Thermus thermophilus SG0.5JP17-16



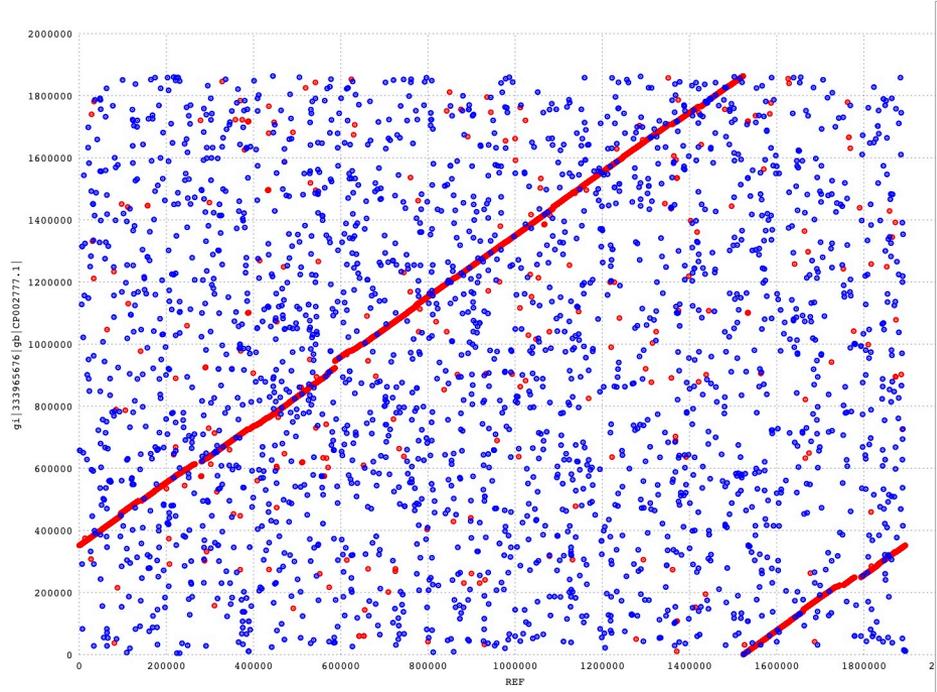
Window=10000 , printed every 100

Usually one plots the Cumulative Strand Bias to more clearly see the turning points

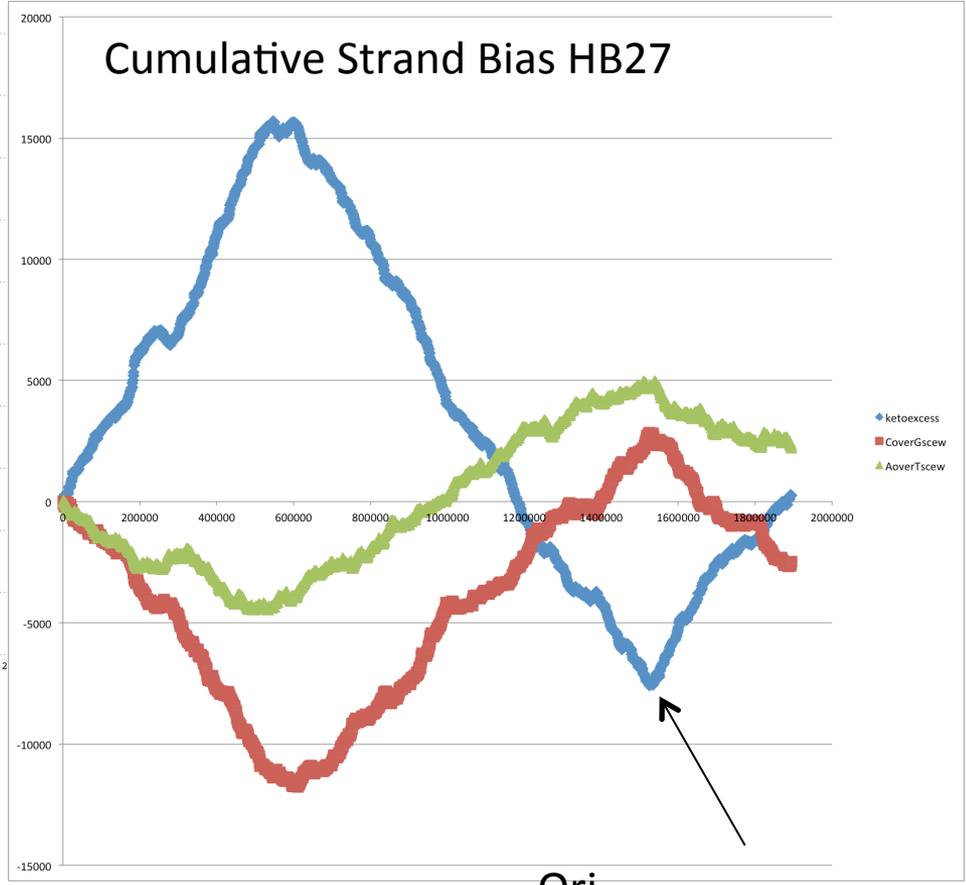


Thermus thermophilus SG0.5JP17-16

Usually, *.fna files of bacterial genomes start with the origin of replication, and the direction is chosen so that the first encoded protein is DnaA (chromosomal replication initiator protein). Sometimes things go wrong.



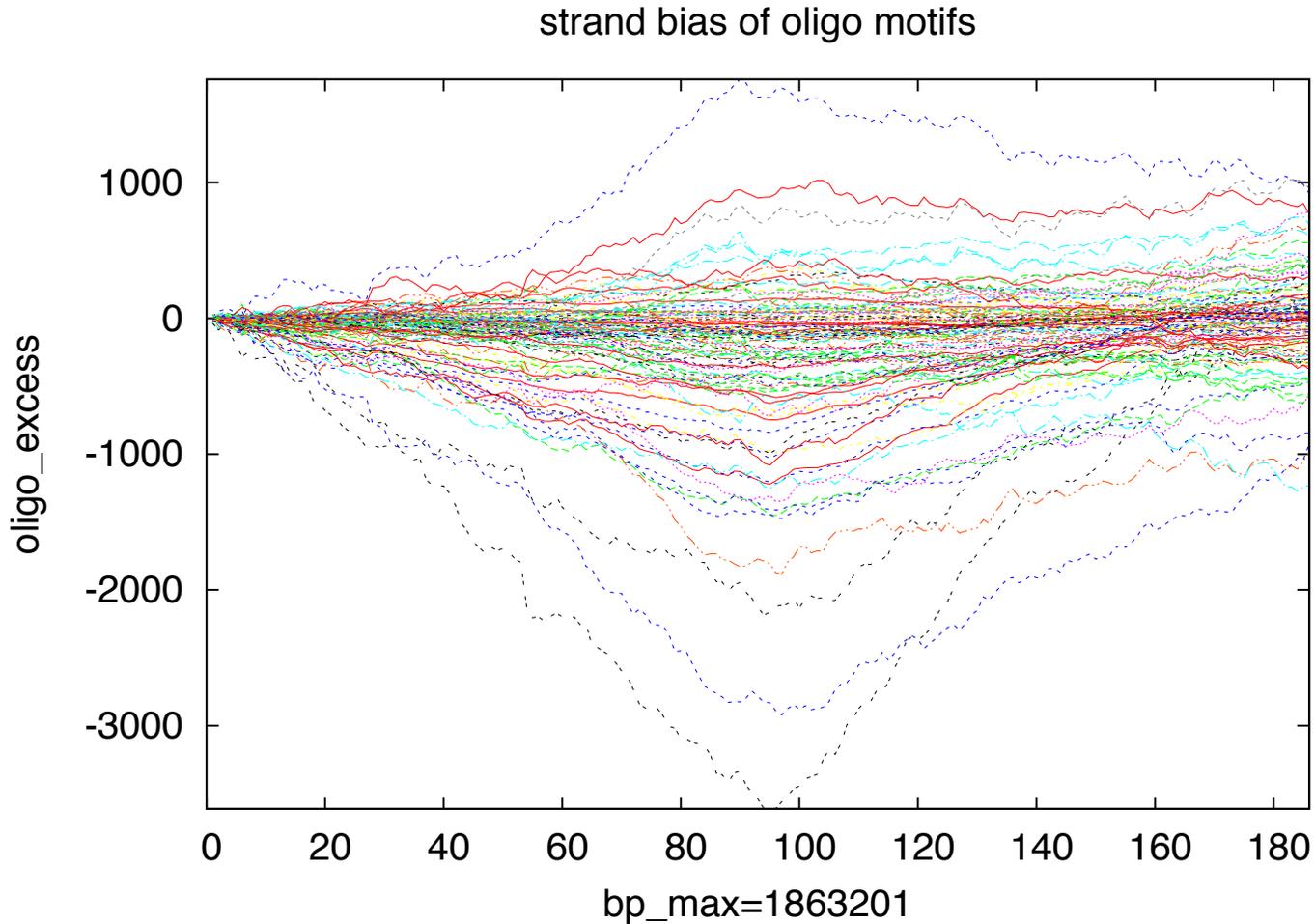
Mummer Plot: HB27 versus SGO



Ori
Should be here

The same can be done with oligonucleotide bias (how often does an oligonucleotide occur on one strand minus occurrence on the other strand)

Tetramer bias for *Thermus thermophilus* SGO



Download a bacterial or archaeal genome of your choice from

<ftp://ftp.ncbi.nih.gov/genomes/Bacteria/>

We are only interested in the sequence of the main chromosome

Write a script that

1. Reads in a genome
2. Assigns the nucleotides of the genome to an array (e.g. @genome)
3. Goes through the @genome array and sequentially counts the numbers of Gs Cs As and Ts
4. Every 1000 (or 5000) nucleotides calculate the excess of Gs over Cs, As over Ts, and the excess of keto bases (G+T-A-C). Feel free to explore other biases.
5. Print the results into a table
6. Plot the columns of the table in Excel or gnuplot

Possibility for

1. Read in a genome

The easiest will be to plagiarize a script you already wrote. The following works, if the script and the *.fna file are in the same directory.

1a: open input and output file, reset stuff

```
#!/usr/bin/perl -w
#initialize genome name and base_hash
$my_genome = "";
%base_hash=();

#assign genome name to $my_genome
@dir=`ls`; # see P24 in the UNIC Perl primer
foreach (@dir) {
    if (m /\.fna$/) {if ($my_genome) {die "More than one genome in directory"} else {$my_genome=($_)}}
}
#####
chomp ($my_genome);
print "\n\n$my_genome is the file name of the genome to be analyzed \n";

# open my genome for input
open (IN, "< $my_genome") or die "cannot open $my_genome:$!";

# open my_table for output
open (OUT, ">my_table" ) or die "cannot open my_table" ;
print OUT "number \tketoexcess\tCoverGscew\tAoverTscew\tbase_hash{A}\tbase_hash{T}\tbase_hash{G}\tbase_hash{C}\n";
# if we want to use exel, we can print a header in the first line";
# if we use gnuplot, we want to omit the header
```

Possibility for

1b: read genome into array

You have 2 possibilities either read and analyze the genome line by line, or read in everything and then start the analysis.

```
$header = <IN>;  
#reads first line of file, next command test for fasta commentline  
if ($header =~m/^>/) {print "\nthe analyzed genome has the following comment line:\n$header \n\n"};  
if (!$header =~m/^>/) {print "this is not in FASTA format \n\n";  
    exit;}  
###      exit - could have died instead;  
  
### Read genome into array @genome  
$number=0;  
  
while (defined ($line=<IN>)){  
  
    #initialise @bases within loop  
  
    @bases=();  
    chomp($line);  
  
    @bases=split(//,$line);  
  
    foreach (@bases) {  
        $number += 1;  
        $genome[$number]=$_  
    }  
}
```

The \$_ is the variable the perl goes through in the foreach loop, see P13

Possibility for

1b': read genome into array

You have 2 possibilities either read and analyze the genome line by line, or read in everything and then start the analysis.

```
if ($header =~m/^>/) {print "\nthe analyzed genome has the following comment line:\n$header \n\n"};

if (!(($header =~m/^>/)) {print "this is not in FASTA format \n\n";
    exit;}}
###      exit - could have died instead;

$number=0;

while (defined ($line=<IN>)){

#initialise @bases within loop
# potential problem: this reads and analyses line by line.
#It might be better, especially if one wants to use nucleotide pairs or oligod, to read everthing in first
    @bases=();
    chomp($line);

    @bases=split(//,$line);

    foreach (@bases) {

        ....
    }
}
```

Write a script that

1. Reads in a genome
2. Assigns the nucleotides of the genome to an array (e.g. @bases)
3. Goes through the @bases array and sequentially counts the numbers of Gs Cs As and Ts

```
@bases=();  
chomp($line);  
  
@bases=split(//,$line);  
  
foreach (@bases) {  
    $number += 1; # counts position in genome  
    $base_hash{$_} += 1;#counts As,Gs and Cs and Ts  
  
    print OUT "$number\t$ketoexcess\t$gcscew\t$satscew\t$base_hash{A}\t$base_hash{T}\t$base_hash{G}\t$base_hash{C}\n";  
    # if we want to use awk we can print a header in the first line
```

4. Every 1000 nucleotides calculates the excess of Gs over Cs, As over Ts, and the excess of keto bases (G+T-A-C). Feel free to explore other biases.

```
#every 1000 nucleotides, print stuff to file  
if ($number%1000==0){  
    $gcscew=($base_hash{C}-$base_hash{G});  
    $satscew=($base_hash{A}-$base_hash{T});  
    $ketoexcess=($base_hash{G}+$base_hash{T})-($base_hash{A}+$base_hash{C});
```

5. Print the results into a table

```
print OUT "$number\t$ketoexcess\t$gcscew\t$satscew\t$base_hash{A}\t$base_hash{T}\t$base_hash{G}\t$base_hash{C}\n";  
# if we want to use awk we can print a header in the first line
```

Close loop(s), close files

6. Plot the columns of the table in Excel or gnuplot

Links to info on gnuplot is [here](#)

If gnuplot is installed on your computer and able to communicate with your x11 terminal program, then

```
> gnuplot
```

Will invoke the gnuplot program

```
> plot "my_table" using 1:2 with lines
```

will plot the 2nd columns against the 1st

```
>set terminal x11
```

Will set the output to x11 (screen)

```
>set terminal png
```

Will set the output to a png file

```
>set terminal postscript
```

sets the output to a postscript file

```
> set out "myplot.png"
```

directs output to the file myplot.png

```
> set multiplot
```

Plots multiple curves into the same figure

```
> plot "my_table" using 1:4 with lines,\
```

```
> "my_table" using 1:2 with lines,\
```

```
> "my_table" using 1:3 with lines
```

plots multiple curves onto the same figure.

gnuplot is installed on the cluster, but you need to direct the output to a file, which is inconvenient (extra credit, if you can make the cluster output to an x11-window on your laptop)

Often one uses a perl script to do the plotting example [here](#).