# Find-A-Gene Project Assignment

Name:     Steven Gan              Quarter:      Fall
PID:      A59020397              Course:       BGGN-213
E-mail:   digan@ucsd.edu        Instructor:   Dr. Barry Grant

## Q1:

Name:        CCCTC-binding factor (CTCF)

Accession:   isoform 1          isoform 2            isoform 3
             NP_006556.1        NP_001177951.1       NP_001350845.1

Species:     *Homo Sapiens*

Function:    DNA insulation; RNA binding; RNA splicing; DNA loop extrusion;
             transcriptional regulation; genome instability.

## Q2:

Isoform 1 will be used for downstream analysis, as it is the longest isoform.

Query:       NP_006556.1

Method:      TBLASTN (2.13.0+) search against with default parameters

Database:    Expressed Sequence Tags (est)

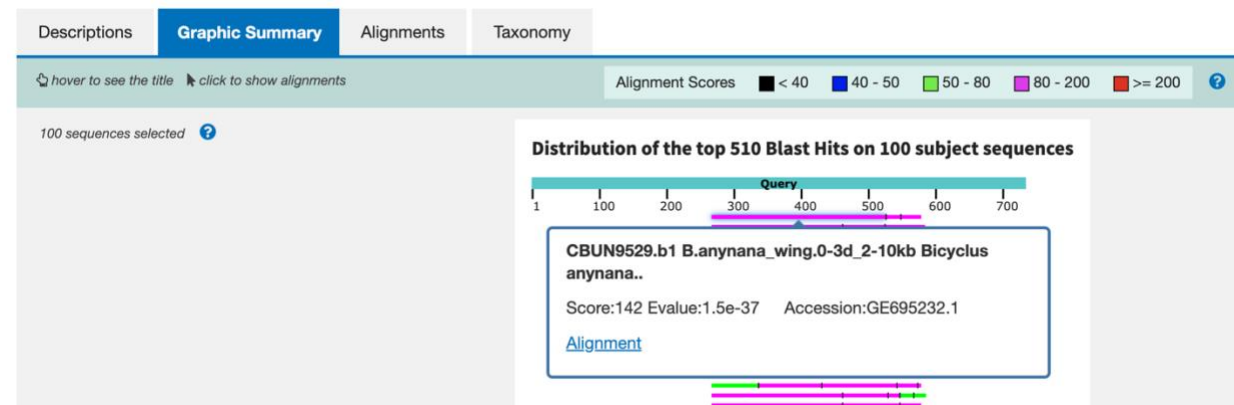Organism:    Papilionoidea (taxid:37572)

TBLASTN Setting:



**Chosen Match:** Accession GE695232.1, a 750 base pair clone from *Bicyclus anynana*. Alignment details see below.

| Descriptions | Graphic Summary | **Alignments** | Taxonomy |

Alignment view [ Pairwise ⌄ ]  ❓ [Restore defaults]     Download ⌄

100 sequences selected ❓

⬇ Download ⌄     GenBank  Graphics    Sort by: [ E value ⌄ ]     ▼ Next ▲ Previous ◄ Descriptions

**CBUN9529.b1 B.anynana_wing.0-3d_2-10kb Bicyclus anynana cDNA clone CBUN9529, mRNA sequence**

Sequence ID: GE695232.1  Length: 750  Number of Matches: 3

Range 1: 3 to 734 GenBank  Graphics     ▼ Next Match ▲ Previous Match

| Score | Expect | Method | Identities | Positives | Gaps | Frame |
|---|---|---|---|---|---|---|
| 154 bits(389) | 1e-41 | Compositional matrix adjust. | 95/284(33%) | 144/284(50%) | 47/284(16%) | +3 |

```
Query  267  QCELCSYTCPRRSNLDRHMKSHTDERPHKCHL----CGRAFRTVTLLRNHLNTHTGTRPH  322
            C++C Y C +R NL  H+++HT E+P+ C +    C R R    LR+H+ THTG +P
Sbjct  3    SCDICHYKCAQRGNLVCHIRTHTCEKPYSCEMYNYKCARKSR----LRHHMTTHTGEKPF  170

Query  323  KCPDCDMAFVTSGELVRHRRYKHTHEKPFKCSMCDYASVEVSKLKRHIRSHTGERPFQCS  382
            C  C+       LV H R  HT EKPF C +C+Y     L  H+++HTGE+PF C
Sbjct  171  SCGICNYKTGVKNSLVCHLR-THTGEKPFCCEICNYKFALKRNLLNHMKTHTGEKPFSCG  347

Query  383  LCSYASRDTYKLKRHMRTHSGEKPYECYICHARFTQSGTMKMHILQKHTENVAKFHCPHC  442
            +C+Y +     L  H+RTH+GEKP+ C IC+ +F
Sbjct  348  ICNYKTGVKNSLVCHLRTHTGEKPFCCEICNYKF-------------------------  449

Query  443  DTVIARKSDLGVHLRKQHSYIEQGKK---CRYCDAVFHERYALIQHQKSHKNEKRFKCDQ  499
            A K +L  H++    I  G+K  C  C+    + +L+ H ++H  EK F C+
Sbjct  450  ----ALKHNLVNHMK-----IHTGEKPFSCEICNYKTRVKNSLVSHLRTHTGEKPFSCEI  602

Query  500  CDYACRQERHMIMHKRTHTGEKPYACSHCDKTFRQKQLLDMHFK  543
            C+Y   ++R+++ H +THTGEKP++C  C+      K  L  H +
Sbjct  603  CNYKSARKRYLLNHMKTHTGEKPFSCDICNYKTGIKNSLVRHMR  734
```

Range 2: 84 to 749 GenBank  Graphics     ▼ Next Match ▲ Previous Match ⚓ First Match

| Score | Expect | Method | Identities | Positives | Gaps | Frame |
|---|---|---|---|---|---|---|
| 142 bits(358) | 1e-37 | Compositional matrix adjust. | 89/255(35%) | 125/255(49%) | 33/255(12%) | +3 |

```
Query  266  FQCELCSYTCPRRSNLDRHMKSHTDERPHKCHLCGRAFRTVTLLRNHLNTHTGTRPHKCP  325
            + CE+ +Y C R+S L  HM +HT E+P  C +C      L  HL THTG +P C
Sbjct  84   YSCEMYNYKCARKSRLRHHMTTHTGEKPFSCGICNYKTGVKNSLVCHLRTHTGEKPFCCE  263

Query  326  DCDMAFVTSGELVRHRRYKHTHEKPFKCSMCDYASVEVSKLKRHIRSHTGERPFQCSLCS  385
            C+  F     L+ H +  HT EKPF C +C+Y +     + L  H+R+HTGE +PF C +C+
Sbjct  264  ICNYKFALKRNLLNHMK-THTGEKPFSCGICNYKTGVKNSLVCHLRTHTGEKPFCCEICN  440

Query  386  YASRDTYKLKRHMRTHSGEKPYECYICHARFTQSGTMKMHILQKHTENVAKFHCPHCDTV  445
            Y     + L  HM+ H+GEKP+ C IC+ +     ++  H L+ HT    F C  C+
Sbjct  441  YKFALKHNLVNHMKIHTGEKPFSCEICNYKTRVKNSLVSH-LRTHTGE-KPFSCEICNYK  614

Query  446  IARKSDLGVHLRKQHSYIEQGKKCRYCDAVFHERYALIQHQKSHKNEKRFKCDQCDYACR  505
            ARK             RY         L+ H K+H  EK F CD C+Y
Sbjct  615  SARK--------------------RY----------LLNHMKTHTGEKPFSCDICNYKTG  704

Query  506  QERHMIMHKRTHTGE  520
            +  ++ H R HTGE
Sbjct  705  IKNSLVRHMRIHTGE  749
```

Range 3: 6 to 728 GenBank  Graphics     ▽ Next Match ▲ Previous Match ⚓ First Match

| Score | Expect | Method | Identities | Positives | Gaps | Frame |
|---|---|---|---|---|---|---|
| 136 bits(343) | 2e-35 | Compositional matrix adjust. | 87/250(35%) | 126/250(50%) | 9/250(3%) | +3 |

```
Query  324  CPDCDMAFVTSGELVRHRRYKHTHEKPFKCSMCDYASVEVSKLKRHIRSHTGERPFQCSL  383
            C C        G LV H R  HT EKP+ C M +Y       S+L+ H+ +HTGE+PF C +
Sbjct  6    CDICHYKCAQRGNLVCHIR-THTCEKPYSCEMYNYKCARKSRLRHHMTTHTGEKPFSCGI  182

Query  384  CSYASRDTYKLKRHMRTHSGEKPYECYICHARFTQSGTMKMHILQKHTENVAKFHCPHCD  443
            C+Y +     L  H+RTH+GEKP+ C IC+ +F    +  H ++ HT TGE-KPFSCGICN
Sbjct  183  CNYKTGVKNSLVCHLRTHTGEKPFCCEICNYKFALKRNLLNH-MKTHTGE-KPFSCGICN  356

Query  444  TVIARKSDLGVHLRKQHSYIEQGKKCRYCDAVFHERYALIQHQKSHKNEKRFKCDQCDYA  503
            K+ L  HLR     E+  C  C+  F  ++ L+ H K H   EK F+C   C+Y
Sbjct  357  YKTGVKNSLVCHLRTHTG--EKPFCCEICNYKFALKHNLVNHMKIHTGEKPFSCEICNYK  530

Query  504  CRQERHMIMHKRTHTGEKPYACSHCDKTFRQKQLLDMHFKRYHDPNFVPAAFVCSKCGKT  563
            R +  ++ H RTHTGEKP++C  C+      +K+ L  H K +          F C  C
Sbjct  531  TRVKNSLVSHLRTHTGEKPFSCEICNYKSARKRYLLNHMKTHTG----EKPFSCDICNYK  698

Query  564  FTRRNTMARH  573
            +N++ RH
Sbjct  699  TGIKNSLVRH  728
```

Alignment Details:

## CBUN9529.b1 B.anynana_wing.0-3d_2-10kb Bicyclus anynana cDNA clone CBUN9529, mRNA sequence

**Sequence ID: GE695232.1      Length: 750      Number of Matches: 3**

Range 1: 3 to 734

**Alignment statistics for match #1**

| Score | Expect | Method | Identities | Positives | Gaps | Frame |
|---|---|---|---|---|---|---|
| 154 bits (389) | 1e-41 | Compositional matrix adjust. | 95/284(33%) | 144/284(50%) | 47/284(16%) | +3 |

```
Query  267  QCELCSYTCPRRSNLDRHMKSHTDERPHKCHL----CGRAFRTVTLLRNHLNTHTGTRPH  322
            C++C Y C +R NL  H+++HT E+P+ C +     C R  R    LR+H+ THTG +P
Sbjct  3    SCDICHYKCAQRGNLVCHIRTHTCEKPYSCEMYNYKCARKSR----LRHHMTTHTGEKPF  170

Query  323  KCPDCDMAFVTSGELVRHRRYKHTHEKPFKCSMCDYASVEVSKLKRHIRSHTGERPFQCS  382
            C  C+        LV H R  HT EKPF C +C+Y        L  H+++HTGE+PF C
Sbjct  171  SCGICNYKTGVKNSLVCHLR-THTGEKPFCCEICNYKFALKRNLLNHMKTHTGEKPFSCG  347

Query  383  LCSYASRDTYKLKRHMRTHSGEKPYECYICHARFTQSGTMKMHILQKHTENVAKFHCPHC  442
            +C+Y +    L  H+RTH+GEKP+ C IC+ +F
Sbjct  348  ICNYKTGVKNSLVCHLRTHTGEKPFCCEICNYKF-------------------------  449

Query  443  DTVIARKSDLGVHLRKQHSYIEQGKK---CRYCDAVFHERYALIQHQKSHKNEKRFKCDQ  499
              A K +L  H++    I  G+K   C  C+     + +L+ H ++H  EK F C+
Sbjct  450  ----ALKHNLVNHMK-----IHTGEKPFSCEICNYKTRVKNSLVSHLRTHTGEKPFSCEI  602

Query  500  CDYACRQERHMIMHKRTHTGEKPYACSHCDKTFRQKQLLDMHFK  543
            C+Y   ++R+++ H +THTGEKP++C  C+     K  L  H +
Sbjct  603  CNYKSARKRYLLNHMKTHTGEKPFSCDICNYKTGIKNSLVRHMR  734
```

Range 2: 84 to 749

**Alignment statistics for match #2**

| Score | Expect | Method | Identities | Positives | Gaps | Frame |
|---|---|---|---|---|---|---|
| 142 bits (358) | 1e-37 | Compositional matrix adjust. | 89/255(35%) | 125/255(49%) | 33/255(12%) | +3 |

```
Query  266  FQCELCSYTCPRRSNLDRHMKSHTDERPHKCHLCGRAFRTVTLLRNHLNTHTGTRPHKCP  325
            + CE+ +Y C R+S L  HM +HT E+P  C +C           L  HL THTG +P  C
Sbjct  84   YSCEMYNYKCARKSRLRHHMTTHTGEKPFSCGICNYKTGVKNSLVCHLRTHTGEKPFCCE  263

Query  326  DCDMAFVTSGELVRHRRYKHTHEKPFKCSMCDYASVEVSKLKRHIRSHTGERPFQCSLCS  385
             C+  F    L+ H +  HT EKPF C +C+Y +    + L  H+R+HTGE+PF C +C+
Sbjct  264  ICNYKFALKRNLLNHMK-THTGEKPFSCGICNYKTGVKNSLVCHLRTHTGEKPFCCEICN  440

Query  386  YASRDTYKLKRHMRTHSGEKPYECYICHARFTQSGTMKMHILQKHTENVAKFHCPHCDTV  445
```

```
        Y     + L  HM+ H+GEKP+ C IC+ +    ++  H L+ HT     F C  C+
Sbjct  441  YKFALKHNLVNHMKIHTGEKPFSCEICNYKTRVKNSLVSH-LRTHTGE-KPFSCEICNYK  614

Query  446  IARKSDLGVHLRKQHSYIEQGKKCRYCDAVFHERYALIQHQKSHKNEKRFKCDQCDYACR  505
             ARK                    RY         L+ H K+H  EK F CD C+Y
Sbjct  615  SARK-------------------RY----------LLNHMKTHTGEKPFSCDICNYKTG  704

Query  506  QERHMIMHKRTHTGE  520
             +  ++ H R HTGE
Sbjct  705  IKNSLVRHMRIHTGE  749
```

Range 3: 6 to 728

**Alignment statistics for match #3**

| Score | Expect | Method | Identities | Positives | Gaps | Frame |
|---|---|---|---|---|---|---|
| 136 bits (343) | 2e-35 | Compositional matrix adjust. | 87/250(35%) | 126/250(50%) | 9/250(3%) | +3 |

```
Query  324  CPDCDMAFVTSGELVRHRRYKHTHEKPFKCSMCDYASVEVSKLKRHIRSHTGERPFQCSL  383
            C C      G LV H R  HT EKP+ C M +Y    S+L+ H+ +HTGE+PF C +
Sbjct  6    CDICHYKCAQRGNLVCHIR-THTCEKPYSCEMYNYKCARKSRLRHHMTTHTGEKPFSCGI  182

Query  384  CSYASRDTYKLKRHMRTHSGEKPYECYICHARFTQSGTMKMHILQKHTENVAKFHCPHCD  443
            C+Y +     L  H+RTH+GEKP+ C IC+ +F    + H ++ HT     F C  C+
Sbjct  183  CNYKTGVKNSLVCHLRTHTGEKPFCCEICNYKFALKRNLLNH-MKTHTGE-KPFSCGICN  356

Query  444  TVIARKSDLGVHLRKQHSYIEQGKKCRYCDAVFHERYALIQHQKSHKNEKRFKCDQCDYA  503
               K+ L  HLR    E+   C  C+  F  ++ L+ H K H  EK F C+ C+Y
Sbjct  357  YKTGVKNSLVCHLRTHTG--EKPFCCEICNYKFALKHNLVNHMKIHTGEKPFSCEICNYK  530

Query  504  CRQERHMIMHKRTHTGEKPYACSHCDKTFRQKQLLDMHFKRYHDPNFVPAAFVCSKCGKT  563
             R +  ++ H RTHTGEKP++C  C+     +K+ L  H K +        F C  C
Sbjct  531  TRVKNSLVSHLRTHTGEKPFSCEICNYKSARKRYLLNHMKTHTG----EKPFSCDICNYK  698

Query  564  FTRRNTMARH  573
             +N++ RH
Sbjct  699  TGIKNSLVRH  728
```

# Q3:

**Chosen sequence:**

```
>B. anynana protein (from BLAST results)
SCDICHYKCAQRGNLVCHIRTHTCEKPYSCEMYNYKCARKSRLRHHMTTHTGEKPFSCGICNYKTGVKNSLVCHLRT
HTGEKPFCCEICNYKFALKRNLLNHMKTHTGEKPFSCGICNYKTGVKNSLVCHLRTHTGEKPFCCEICNYKFALKHN
LVNHMKIHTGEKPFSCEICNYKTRVKNSLVSHLRTHTGEKPFSCEICNYKSARKRYLLNHMKTHTGEKPFSCDICNY
KTGIKNSLVRHMRIHTG
```

## All six reading frame:

Name:       *Bicyclus* CTCF-like protein

Species:    *Bicyclus anynana*

Taxonomy:   Eukaryota; Metazoa; Ecdysozoa; Arthropoda; Hexapoda; Insecta;
            Pterygota; Neoptera; Endopterygota; Lepidoptera; Glossata;
            Ditrysia; Papilionoidea; Nymphalidae; Satyrinae; Satyrini;
            Mycalesina; Bicyclus.

# Q4:

BLASTP search on non-redundant protein sequences (nr) hits top on zinc finger protein 84-like proteins on *Bicyclus anynana*, with identity percentage of 95.16%, suggesting a possible novel protein. See details below.

BLASTP setting:

## BLASTP results:

| Descriptions | Graphic Summary | Alignments | Taxonomy |
|---|---|---|---|

**Sequences producing significant alignments**    Download ⌄    Select columns ⌄    Show 100 ▼    ❓

☑ select all  *100 sequences selected*    GenPept  Graphics    Distance tree of results    Multiple alignment  MSA Viewer

| Description | Scientific Name | Max Score | Total Score | Query Cover | E value | Per. Ident | Acc. Len | Accession |
|---|---|---|---|---|---|---|---|---|
| ☑ zinc finger protein 84-like [Bicyclus anynana] | Bicyclus anynana | 488 | 1550 | 100% | 1e-165 | 95.16% | 810 | XP_023937889.1 |
| ☑ histone-lysine N-methyltransferase PRDM9-like isoform X1 [Maniola jurtina] | Maniola jurtina | 337 | 1414 | 100% | 4e-113 | 62.50% | 301 | XP_045778716.1 |
| ☑ zinc finger protein 260-like isoform X1 [Bicyclus anynana] | Bicyclus anynana | 334 | 1907 | 100% | 9e-107 | 65.06% | 704 | XP_023953410.1 |
| ☑ zinc finger protein 260-like isoform X2 [Bicyclus anynana] | Bicyclus anynana | 334 | 1906 | 100% | 1e-106 | 65.06% | 701 | XP_023953411.1 |
| ☑ gastrula zinc finger protein XlCGF57.1-like [Bicyclus anynana] | Bicyclus anynana | 318 | 1852 | 100% | 3e-104 | 62.75% | 411 | XP_023952909.1 |
| ☑ gastrula zinc finger protein XlCGF8.2DB-like isoform X1 [Nilaparvata lugens] | Nilaparvata lugens | 316 | 887 | 100% | 1e-103 | 55.69% | 381 | XP_039298152.1 |
| ☑ zinc finger protein 182-like isoform X4 [Bicyclus anynana] | Bicyclus anynana | 325 | 611 | 100% | 2e-103 | 63.82% | 684 | XP_023952407.1 |
| ☑ zinc finger protein 182-like isoform X3 [Bicyclus anynana] | Bicyclus anynana | 325 | 612 | 100% | 4e-103 | 63.82% | 729 | XP_023952406.1 |
| ☑ zinc finger protein 182-like isoform X1 [Bicyclus anynana] | Bicyclus anynana | 325 | 612 | 100% | 4e-103 | 63.82% | 732 | XP_023952404.1 |
| ☑ zinc finger protein 260-like isoform X2 [Bicyclus anynana] | Bicyclus anynana | 325 | 612 | 100% | 4e-103 | 63.82% | 732 | XP_023952405.1 |
| ☑ gastrula zinc finger protein XlCGF57.1-like isoform X2 [Maniola jurtina] | Maniola jurtina | 307 | 756 | 100% | 2e-102 | 62.61% | 245 | XP_045778717.1 |
| ☑ gastrula zinc finger protein XlCGF57.1-like [Nilaparvata lugens] | Nilaparvata lugens | 319 | 1782 | 100% | 4e-102 | 56.68% | 619 | XP_022189878.2 |
| ☑ gastrula zinc finger protein XlCGF17.1-like [Maniola hyperantus] | Maniola hyperantus | 309 | 1708 | 100% | 8e-102 | 59.27% | 329 | XP_034839155.1 |
| ☑ gastrula zinc finger protein XlCGF8.2DB-like isoform X1 [Bicyclus anynana] | Bicyclus anynana | 310 | 610 | 100% | 9e-102 | 62.55% | 364 | XP_023950901.1 |

## Alignment details:

| Descriptions | Graphic Summary | **Alignments** | Taxonomy |
|---|---|---|---|

Alignment view  [ Pairwise ▼ ]   ❓ [Restore defaults]    Download ⌄

*100 sequences selected* ❓

⬇ Download ⌄   GenPept  Graphics   Sort by: [ E value ▼ ]    ▼ Next  ▲ Previous  ◄Descriptions

**zinc finger protein 84-like [Bicyclus anynana]**
Sequence ID: XP_023937889.1  Length: **810**  Number of Matches: **4**

Range 1: 476 to 723 GenPept  Graphics    ▼ Next Match  ▲ Previous Match

**Related Information**
Gene - associated gene details
Genome Data Viewer - aligned genomic context

| Score | Expect | Method | Identities | Positives | Gaps |
|---|---|---|---|---|---|
| 488 bits(1257) | 1e-165 | Compositional matrix adjust. | 236/248(95%) | 238/248(95%) | 0/248(0%) |

```
Query  1    SCDICHYKCAQRGNLVCHIRTHTCEKPYSCEMYNYKCARKSRLRHHMTTHTGEKPFSCGI   60
            SCDICHYKCAQRGNLVCHIRTHT EKPYSCEM NYKCA KSRLRHHMTTHTGEKPFSCGI
Sbjct  476  SCDICHYKCAQRGNLVCHIRTHTGEKPYSCEMCNYKCAHKSRLRHHMTTHTGEKPFSCGI   535

Query  61   CNYKTGVKNSLVCHLRTHTGEKPFCCEICNYKFALKRNLLNHMKTHTGEKPFSCGICNYK   120
            CNYKTGVKNSLVCHLRTHTGEKPFCCEICNYKFALKRNLLNHMKTHTGEKPFSCGICNYK
Sbjct  536  CNYKTGVKNSLVCHLRTHTGEKPFCCEICNYKFALKRNLLNHMKTHTGEKPFSCGICNYK   595

Query  121  TGVKNSLVCHLRTHTGEKPFCCEICNYKFALKHNLVNHMKIHTGEKPFSCEICNYKTRVK   180
            TGVKNSLVCHLRTHTGEKPFCCEICNYKFALKHNLVNHMKIHTGEKPFSCEICNYKTRVK
Sbjct  596  TGVKNSLVCHLRTHTGEKPFCCEICNYKFALKHNLVNHMKIHTGEKPFSCEICNYKTRVK   655

Query  181  NSLVSHLRTHTGEKPFSCEICNYKSARKRYLLNHMKTHTGEKPFSCDICNYKTGIKNSLV   240
            NSLV HLRTHTGEKPF CEICNYK A K  L+NHMKTHTGEKPFSCDICNYKTGIKN+LV
Sbjct  656  NSLVCHLRTHTGEKPFCCEICNYKFALKHNLVNHMKTHTGEKPFSCDICNYKTGIKNTLV   715

Query  241  RHMRIHTG   248
            RHMR HTG
Sbjct  716  RHMRTHTG   723
```

**Range 2: 532 to 779** GenPept Graphics ▼ Next Match ▲ Previous Match ⚓ First Match

| Score | Expect | Method | Identities | Positives | Gaps |
|---|---|---|---|---|---|
| 409 bits(1051) | 1e-134 | Compositional matrix adjust. | 200/248(81%) | 216/248(87%) | 0/248(0%) |

```
Query  1    SCDICHYKCAQRGNLVCHIRTHTCEKPYSCEMYNYKCARKSRLRHHMTTHTGEKPFSCGI  60
            SC IC+YK   + +LVCH+RTHT EKP+ CE+ NYK A K  L +HM THTGEKPFSCGI
Sbjct  532  SCGICNYKTGVKNSLVCHLRTHTGEKPFCCEICNYKFALKRNLLNHMKTHTGEKPFSCGI  591

Query  61   CNYKTGVKNSLVCHLRTHTGEKPFCCEICNYKFALKRNLLNHMKTHTGEKPFSCGICNYK  120
            CNYKTGVKNSLVCHLRTHTGEKPFCCEICNYKFALK NL+NHMK HTGEKPFSC ICNYK
Sbjct  592  CNYKTGVKNSLVCHLRTHTGEKPFCCEICNYKFALKHNLVNHMKIHTGEKPFSCEICNYK  651

Query  121  TGVKNSLVCHLRTHTGEKPFCCEICNYKFALKHNLVNHMKIHTGEKPFSCEICNYKTRVK  180
            T VKNSLVCHLRTHTGEKPFCCEICNYKFALKHNLVNHMK HTGEKPFSC+ICNYKT +K
Sbjct  652  TRVKNSLVCHLRTHTGEKPFCCEICNYKFALKHNLVNHMKTHTGEKPFSCDICNYKTGIK  711

Query  181  NSLVSHLRTHTGEKPFSCEICNYKSARKRYLLNHMKTHTGEKPFSCDICNYKTGIKNSLV  240
            N+LV H+RTHTGEKPFSCEICN+KSA K  LL+HMKTHTGEKPFSC ICNYK   K  L+
Sbjct  712  NTLVRHMRTHTGEKPFSCEICNHKSALKHSLLSHMKTHTGEKPFSCKICNYKCVRKQHLL  771

Query  241  RHMRIHTG  248
            HM+ HTG
Sbjct  772  GHMKTHTG  779
```

**Range 3: 588 to 809** GenPept Graphics ▼ Next Match ▲ Previous Match ⚓ First Match

| Score | Expect | Method | Identities | Positives | Gaps |
|---|---|---|---|---|---|
| 334 bits(857) | 8e-106 | Compositional matrix adjust. | 165/222(74%) | 187/222(84%) | 0/222(0%) |

```
Query  1    SCDICHYKCAQRGNLVCHIRTHTCEKPYSCEMYNYKCARKSRLRHHMTTHTGEKPFSCGI  60
            SC IC+YK   + +LVCH+RTHT EKP+ CE+ NYK A K  L +HM  HTGEKPFSC I
Sbjct  588  SCGICNYKTGVKNSLVCHLRTHTGEKPFCCEICNYKFALKHNLVNHMKIHTGEKPFSCEI  647

Query  61   CNYKTGVKNSLVCHLRTHTGEKPFCCEICNYKFALKRNLLNHMKTHTGEKPFSCGICNYK  120
            CNYKT VKNSLVCHLRTHTGEKPFCCEICNYKFALK NL+NHMKTHTGEKPFSC ICNYK
Sbjct  648  CNYKTRVKNSLVCHLRTHTGEKPFCCEICNYKFALKHNLVNHMKTHTGEKPFSCDICNYK  707

Query  121  TGVKNSLVCHLRTHTGEKPFCCEICNYKFALKHNLVNHMKIHTGEKPFSCEICNYKTRVK  180
            TG+KN+LV H+RTHTGEKPF CEICN+K ALKH+L++HMK HTGEKPF+SC+ICNYK   K
Sbjct  708  TGIKNTLVRHMRTHTGEKPFSCEICNHKSALKHSLLSHMKTHTGEKPFSCKICNYKCVRK  767

Query  181  NSLVSHLRTHTGEKPFSCEICNYKSARKRYLLNHMKTHTGEK  222
            L+ H++THTGEK F C++CNYK ARKR L +HMKTHTG K
Sbjct  768  QHLLGHMKTHTGEKSFCCKLCNYKCARKRDLESHMKTHTGGK  809
```

**Range 4: 447 to 667** GenPept Graphics ▼ Next Match ▲ Previous Match ⚓ First Match

| Score | Expect | Method | Identities | Positives | Gaps |
|---|---|---|---|---|---|
| 317 bits(813) | 3e-99 | Compositional matrix adjust. | 161/221(73%) | 182/221(82%) | 3/221(1%) |

```
Query  31   EMYN-YKCARKSRLRHH--MTTHTGEKPFSCGICNYKTGVKNSLVCHLRTHTGEKPFCCE  87
            ++Y+ +K  +K+ L  +  T T E +SC IC+YK   + +LVCH+RTHTGEKP+ CE
Sbjct  447  QLYDIFKKPKKTVLDENPRVKTLTNEILYSCDICHYKCAQRGNLVCHIRTHTGEKPYSCE  506

Query  88   ICNYKFALKRNLLNHMKTHTGEKPFSCGICNYKTGVKNSLVCHLRTHTGEKPFCCEICNY  147
            +CNYK A K  L +HM THTGEKPFSCGICNYKTGVKNSLVCHLRTHTGEKPFCCEICNY
Sbjct  507  MCNYKCAHKSRLRHHMMTHTGEKPFSCGICNYKTGVKNSLVCHLRTHTGEKPFCCEICNY  566

Query  148  KFALKHNLVNHMKIHTGEKPFSCEICNYKTRVKNSLVSHLRTHTGEKPFSCEICNYKSAR  207
            KFALK NL+NHMK HTGEKPFSC ICNYKT VKNSLV HLRTHTGEKPF CEICNYK A
Sbjct  567  KFALKRNLLNHMKTHTGEKPFSCGICNYKTGVKNSLVCHLRTHTGEKPFCCEICNYKFAL  626

Query  208  KRYLLNHMKTHTGEKPFSCDICNYKTGIKNSLVRHMRIHTG  248
            K  L+NHMK HTGEKPFSC+ICNYKT +KNSLV H+R HTG
Sbjct  627  KHNLVNHMKIHTGEKPFSCEICNYKTRVKNSLVCHLRTHTG  667
```

# Q5:

## Sequences for multiple alignment:

```
>Bicyclus_anynana (from BLAST results)
SCDICHYKCAQRGNLVCHIRTHTCEKPYSCEMYNYKCARKSRLRHHMTTHTGEKPFSCGICNYKTGVKNSLVCHLRT
HTGEKPFCCEICNYKFALKRNLLNHMKTHTGEKPFSCGICNYKTGVKNSLVCHLRTHTGEKPFCCEICNYKFALKHN
LVNHMKIHTGEKPFSCEICNYKTRVKNSLVSHLRTHTGEKPFSCEICNYKSARKRYLLNHMKTHTGEKPFSCDICNY
KTGIKNSLVRHMRIHTG

>Homo_sapiens ref|NP_006556.1| transcriptional repressor CTCF isoform 1 [Homo
sapiens]
MEGDAVEAIVEESETFIKGKERKTYQRRREGGQEEDACHLPQNQTDGGEVVQDVNSSVQMVMMEQLDPTLLQMKTEV
MEGTVAPEAEAAVDDTQIITLQVVNMEEQPINIGELQLVQVPVPVTVPVATTSVEELQGAYENEVSKEGLAESEPMI
```

CHTLPLPEGFQVVKVGANGEVETLEQGELPPQEDPSWQKDPDYQPPAKKTKKTKKSKLRYTEEGKDVDVSVYDFEEE
QQEGLLSEVNAEKVVGNMKPPKPTKIKKKGVKKTFQCELCSYTCPRRSNLDRHMKSHTDERPHKCHLCGRAFRTVTL
LRNHLNTHTGTRPHKCPDCDMAFVTSGELVRHRRYKHTHEKPFKCSMCDYASVEVSKLKRHIRSHTGERPFQCSLCS
YASRDTYKLKRHMRTHSGEKPYECYICHARFTQSGTMKMHILQKHTENVAKFHCPHCDTVIARKSDLGVHLRKQHSY
IEQGKKCRYCDAVFHERYALIQHQKSHKNEKRFKCDQCDYACRQERHMIMHKRTHTGEKPYACSHCDKTFRQKQLLD
MHFKRYHDPNFVPAAFVCSKCGKTFTRRNTMARHADNCAGPDGVEGENGGETKKSKRGRKRKMRSKKEDSSDSENAE
PDLDDNEDEEEPAVEIEPEPEPQPVTPAPPPAKKRRGRPPGRTNQPKQNQPTAIIQVEDQNTGAIENIIVEVKKEPD
AEPAEGEEEEAQPAATDAPNGDLTPEMILSMMDR


>Maniola_jurtina ref| XP_045782119.1 | transcriptional repressor CTCF-like
isoform X2 [Maniola jurtina]
MAGICCVDGCDPTAEDVTYFKFPNSRTLRRKWLDAINNSVKVTLDTAVCSRHFLPNQYEVIRGKKRLKAKVVPSVFD
NITKPTSPQKEKTDSSDGEDSVPLQKVKSVATDNTDTGQSKQSPDHRLEDRQDSEASVRNGVKEDDIVKRKQPDIVS
ITDSESNKDIEDIITHYQIKQIRPLHKPTDRTTPDMSIEIEVPLAMDGEMEIGRNREMEMGRNGEIDDVIDVDEEAE
PVFIEVAVGKGGGVEETTNEDCMMLLESVQCEVDPSCLMFPPEDPGNDPGDDAGNDSDVIDLGEKKEDPVSLLTSSD
EDEVIIEEPKYDMVEVSDETDEDDVPLVRLVDKPSQNKFPKNTKNTDILSETNLTKLLWGRLCEYYCLECRFTSTSN
AELRKHMQEHSTQVIQVCEICSYTTSSKHQYIRHKRKHKEDKRFKCHLCKYSARHNMSLIYHLKSHDNGQFVSDMSV
FKCEKCNFETDYKVSLMKHIRICSSKSKRYSCAKCSYETDRRSDLKRHKARKHNTGKDGDYEPPAWVTRAKKPKCDK


>Nilaparvata_lugens ref|XP_039287301.1| transcriptional repressor CTCF
[Nilaparvata lugens]
MSPPDKVQVQTEIKLEDGVTIVPENVTDIQNYLDTFNKEIQGGEQVVQQVGVVAADEGGSEEGTYYVDQAGQYYYQS
ASCDGQQVMTVVSGLPGASGESGESFVALPASAASSQRDNVGGSAPLLIQAATGGGASASGGGTTVGGATVGGASAE
GGGGATYQTVTIVPSETNPGELSYLLIVQQPGDEDEGEDGQKDKDEDDDHDLTVYDFDDAEDVGTVSGMESGDEDDK
SKIVKFMPKKSQTVTQAHMCNYCNYTSPKRYLLSRHMKSHSEERPHKCSVCERGFKTLASLQNHVNTHTGTKPHRCK
HCDSAFTTSGELVRHVRYKHTHEKPHKCTICDYASVELSKMRNHMRCHTGERPYQCPHCTYASPDTFKLKRHLRIHT
GEKPYECDICHARFTQSNSLKAHKLIHSGQLHTSISPPSLLGVLTENSFNDFEDWNDFHFRVDNFSMSALQFAKFSV
HSKSHEGEKCWRCELCPYASVSQRHLESHMLIHTDQKPYQCDQCDQSFRQKQLLRRHQNLYHNPNYVPPPPREKTHE
CPECQRAFRHKGNLIRHLSVHDPESLAQERQMLLKQGRQRKLQNINGQRVEVIPGDEDDEDEDDELNGQVMAVEGSD
GQQYVVLEVIQLQDDNGQEQAVAVMAADGGLEQAVAALHGAGAEDDEEELDEEDEEEDDVHITPDMVEDHEMMSSLR
HQTRQKTQHDMANCFGFDDDEEEEEEDDIGISLKQSNTKSIHLLRSGLQ


>Mus_musculus ref|NP_001390652.1| transcriptional repressor CTCF isoform 1
[Mus musculus]
MEGEAVEAIVEESETFIKGKERKTYQRRREGGQEEDACHLPQNQTDGGEVVQDVNSSVQMVMMEQLDPTLLQMKTEV
MEGTVAPEAEAAVDDTQIITLQVVNMEEQPINIGELQLVQVPVPVTVPVATTSVEELQGAYENEVSKEGLAESEPMI
CHTLPLPEGFQVVKVGANGEVETLEQGELPPQEDSSWQKDPDYQPPAKKTKKTKKSKLRYTEEGKDVDVSVYDFEEE
QQEGLLSEVNAEKVVGNMKPPKPTKIKKKGVKKTFQCELCSYTCPRRSNLDRHMKSHTDERPHKCHLCGRAFRTVTL
LRNHLNTHTGTRPHKCPDCDMAFVTSGELVRHRRYKHTHEKPFKCSMCDYASVEVSKLKRHIRSHTGERPFQCSLCS
YASRDTYKLKRHMRTHSGEKPYECYICHARFTQSGTMKMHILQKHTENVAKFHCPHCDTVIARKSDLGVHLRKQHSY
IEQGKKCRYCDAVFHERYALIQHQKSHKNEKRFKCDQCDYACRQERHMIMHKRTHTGEKPYACSHCDKTFRQKQLLD
MHFKRYHDPNFVPAAFVCSKCGKTFTRRNTMARHADNCAGPDGVEGENGGETKKSKRGRKRKMRSKKEDSSDSEENA
EPDLDDNEEEEPAVEIEPEPEPQPQPPPPPQPVAPAPPPAKKRRGRPPGRTNQPKQNQPTAIIQVEDQNTGAIENI
IVEVKKEPDAEPAEGEEEEAQAATTDAPNGDLTPEMILSMMDR


>Drosophila_melanogaster ref|NP_648109.1| CTCF [Drosophila melanogaster]
MPRRTKKDEDPEDLQTFLNNFHKEIEGNSDEKVVNTILEAISAEAIDLDENGAEEAGGSKPMEEAEADLDHAEEAEEE
EEDDEDKYFIDDEGNCYIKTTPKKQKELQKKLKQAAAKPGKATRSVVSTATNKSINLRPAKSTPKATTSKPPPEPKA
ISVRPARAAAAKAKQSAMPPPPALVVKVPAPRGRPRKNPVIPKPEPMDLERELEELVDEPDISSMVTELSDYTVDEA
AVEAATATLTPNEAEVYEFEDNATTEDENADKKDVDFVLSNKEVKLKTASSTSQNSNASGHKYSCPHCPYTASKKFL
ITRHSRSHDVEPSFKCSICERSFRSNVGLQNHINTHMGNKPHKCKLCESAFTTSGELVRHTRYKHTKEKPHKCTECT
YASVELTKLRRHMTCHTGERPYQCPHCTYASQDMFKLKRHMVIHTGEKKYQCDICKSRFTQSNSLKAHKLIHSVVDK
PVFQCNYCPTTCGRKADLRVHIKHMHTSDVPMTCRRCGQQLPDRYQYKLHVKSHEGEKCYSCKLCSYASVTQRHLAS
HMLIHLDEKPFHCDQCPQAFRQRQLLRRHMNLVHNEEYQPPEPREKLHKCPSCPREFTHKGNLMRHMETHDDSANAR
EKRRRLKLGRNVRLQKDGTVITLIKDQYVDMDRDQEENEEDDNPESYDLAEIEPENSEAEDADDDVETIVSDPIRQR
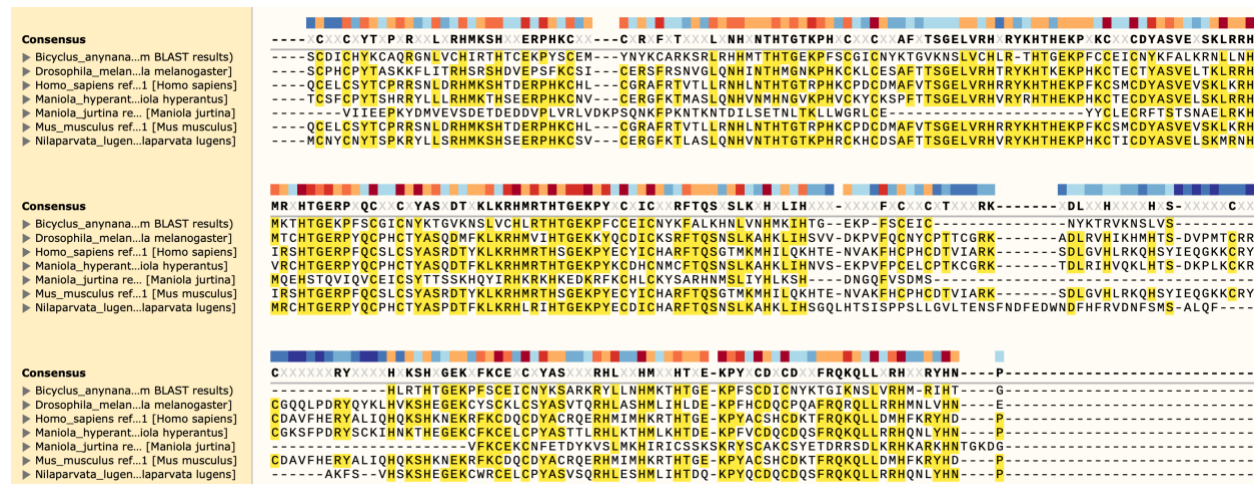IKPAPIIINKQARLAASEKQPMIINQRLRSQRGTKTFHIKEEPDNSDFTVEWGDDGEVMVVELVNGDEEVLVKHEP

SANSKISAKNCFGFEDDDDYEEYGDGENEVDGASQEFLQLMDMIEQDS

>Maniola_hyperantus ref|XP_034839420.1| transcriptional repressor CTCF-like
[Maniola hyperantus]
MPPPDKKSANKCKTILQTYLNSFDQDNEPTTVIVNGDGDEADAGVTYFVDEEGRYYYQPAGDSQNLVSLPIEAEAED
GTEIPQEAQMLVDGDGYQTVTLMPSEEGGELSYVLVMQEETKPVMNIDIKVDQDEEKSSDVYKFEEEEEEDPPIEVS
DEVEESIKPKLTFAMKRSKHLRPSFTCSFCPYTSHRRYLLLRHMKTHSEERPHKCNVCERGFKTMASLQNHVNMHNG
VKPHVCKYCKSPFTTSGELVRHVRYRHTHEKPHKCTECDYASVELSKLRRHVRCHTGERPYQCPHCTYASQDTFKLK
RHMRTHTGEKPYKCDHCNMCFTQSNSLKAHKLIHNVSEKPVFPCELCPTKCGRKTDLRIHVQKLHTSDKPLKCKRCG
KSFPDRYSCKIHNKTHEGEKCFKCELCPYASTTLRHLKTHMLKHTDEKPFVCDQCDQSFRQKQLLRRHQNLYHNPNY
EPKPPKEKTHTCHECKRTFAHKGNLIRHLAIHDPDSGHQERALALRLGRQKKIKFVDGNVKTDDSDNEPEEIMKLDL
GGNQLERGELLTVADNDGQQYVVLEVIQAEDGETQIVSAADYEEEEEEEEEEDEDDEELDKKEIIYEQIKPKGMME
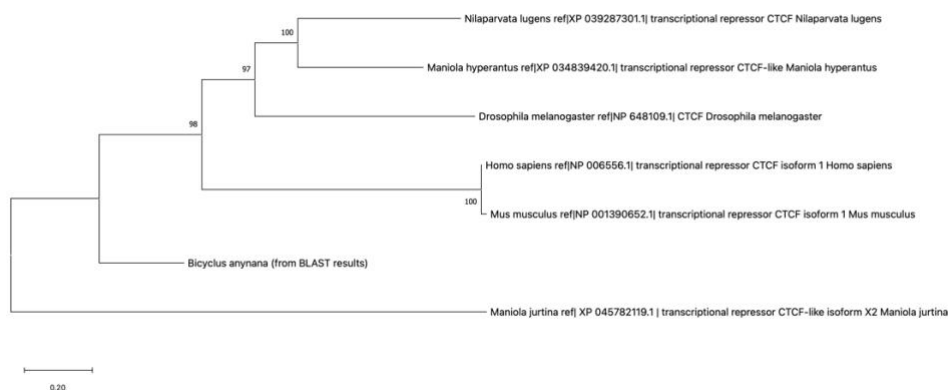RTIKLESDVDTCFGFDEDEEEPDEDEEGIAYNDKIVLRIV

## Alignment:
Obtained using MUSCLE (version 3.8) in SnapGene (version 6.0.2)
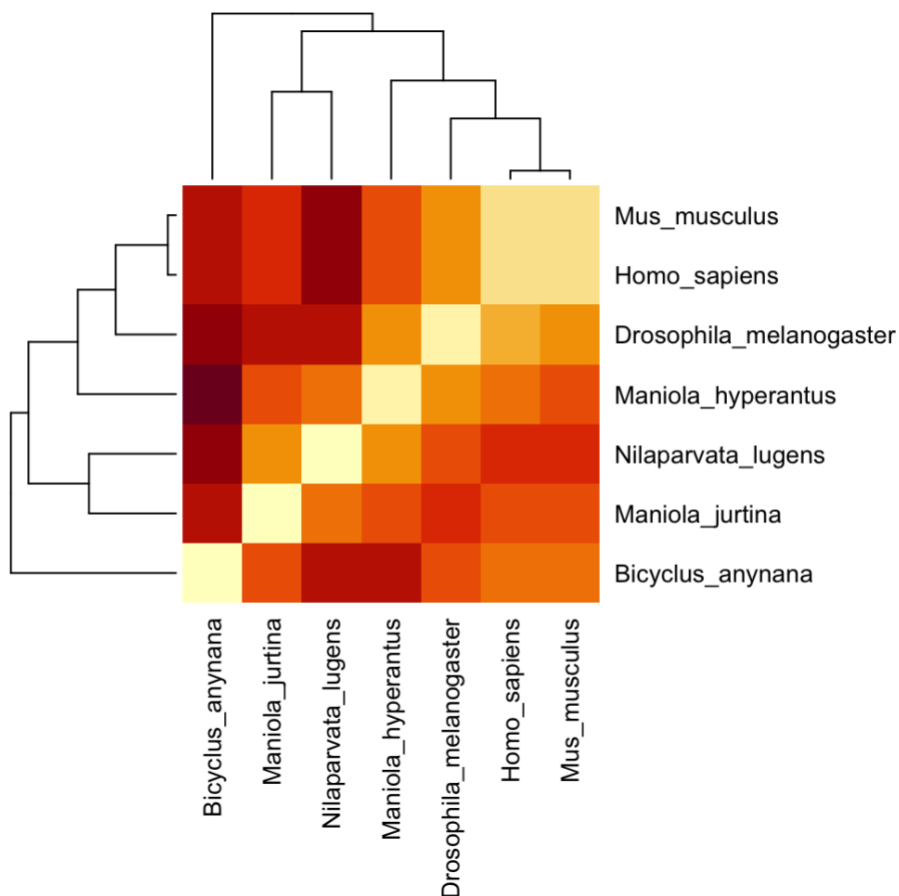


# Q6:

Phylogenetic tree: Generated using MEGA (11.0.13), aligned with MUSCLE (3.8)

# Q7:

Heatmap:



# Q8:

Top three hits of blast search based on the consensus sequence:

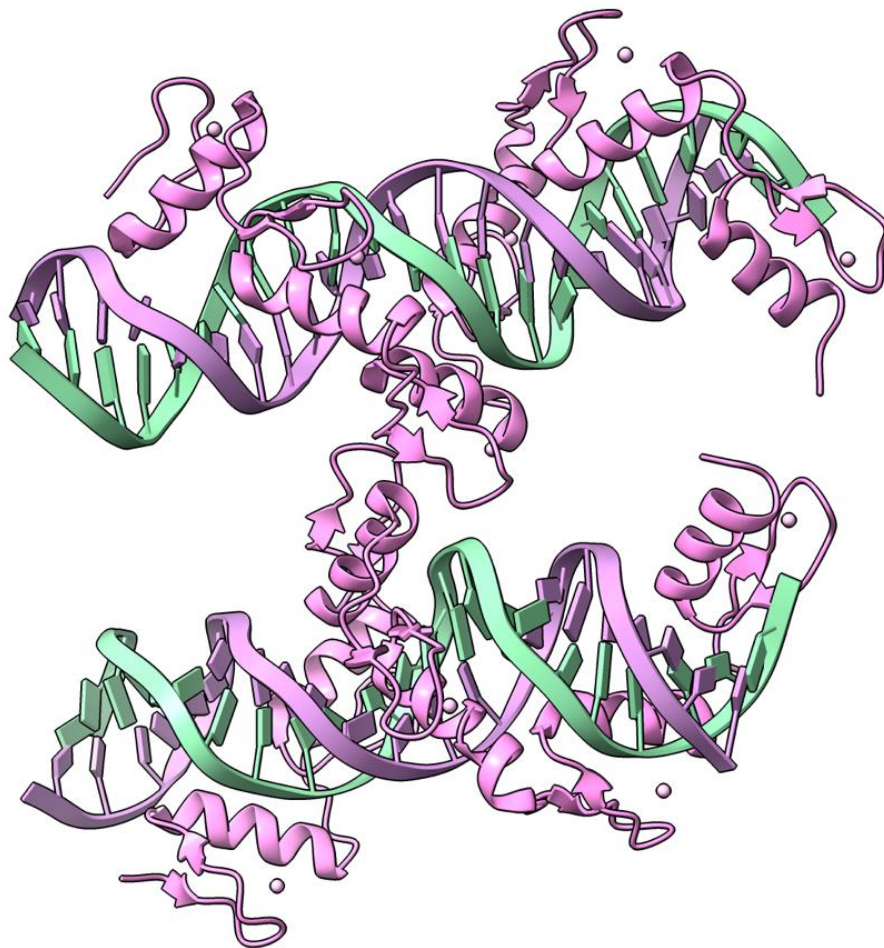| ID | Technique | Resolution | Source | E-value | Identity (%) |
|---|---|---|---|---|---|
| 6QNX_C | X-ray Diffraction | 2.700 | *Homo sapiens* | 1.21e-95 | 56.584 |
| 5YEF_A | X-ray Diffraction | 2.807 | *Homo sapiens* | 5.49e-66 | 57.071 |
| 5T0U_A | X-ray Diffraction | 3.199 | *Homo sapiens* | 1.32e-62 | 61.677 |

Consensus sequence:

```
>Aligned Consensus (threshold of >50%)
XCXXCXYTXPXRXXLXRHMKSHXXERPHKCXXCXRXFXTXXXLXNHXNTHTGTKPHXCXXCXXAFXTSGELVRHXRY
KHTHEKPXKCXXCDYASVEXSKLRRHMRXHTGERPXQCXXCXYASXDTXKLRHMRTHTGEKPYXCXICXXRFTQSX
SLKXHXLIHXXXXXXXFXCXXCXTXXXXRKXDLXXHXXXXHXSXXXXXCXXCXXXXXXRYXXXXHXKSHXGEKXFKCE
XCXYASXXXRHLXXHMXXHTXEKPYXCDXCDXXFRQKQLLXRHXXRYHNP
```

# Q9:

5T0U might not bear too much similarity to *Bicyclus anynana* CTCF like proteins since the identity only go as high as 62%. Figure below is the structure of homodimer CTCF proteins in human binding to two DNA strands.

Structure of 5T0U:

# Q10:

https://www.ebi.ac.uk/chembl/target_report_card/CHEMBL4523233/

Target report of CTCF protein in ChEMBL:

Only one chemical, CHEMBL2334661 ($C_{19}H_{16}O_3$; name undefined), is reported to inhibit CTCF in HUVEC (*Homo Sapiens* cell line) by reducing *CTCF* gene transcription.

Assay ID: CHEMBL4421277

| | |
|---|---|
| **Assay ID:** | CHEMBL4421277 |
| **Type:** | Binding |
| **Description:** | Inhibition of CTCF in HUVEC assessed as reduction of CTCF transcriptional activity by genome-wide RNA-seq and ChIP-seq analysis |
| **Format:** | BAO_0000219 |
| **Journal:** | No Reference Available |
| **Organism:** | Homo sapiens |
| **Strain:** | --- |
| **Tissue:** | --- |
| **Cell Type:** | HUVEC |
| **Subcellular Fraction:** | --- |
| **Target:** | CHEMBL4523233 |
| **Document:** | CHEMBL4420080 |
| **Cell:** | CHEMBL3307501 |
| **Tissue:** | |