

Repeats and homorepeats

Miguel Andrade
Faculty of Biology,
Johannes Gutenberg University
Mainz, Germany
andrade@uni-mainz.de

Repeats

Frequency

14% proteins contains repeats (Marcotte et al, 1999)

1: Single amino acid repeats.

2: Longer imperfect tandem repeats.
Assemble in structure.

Definition repeats

Sequence, long, imperfect, tandem

MRAVVKSPIMCHEKSPSVCSPLNMTSSVCSPAGINSVSSTTASF
GSFPVHSPITQGTPLTCSPNVENRGSRSHSPAASNMGSPPLSSP
LSSMKSSISSLPPSHCSVKSPVSSPNNVTLRSSVSSPANINN

Definition repeats

Sequence, long, imperfect, tandem

MRAVVK**SP**IMCHEKSPSVC**SP**LNMTSSVC**SP**AGINSVSSTTASF
GSFPVH**SP**ITQGTPLTC**SP**NVENRGSRSH**SP**AHASNVGSPLS**SP**
LSSMKSSIS**SP**PSHCSVKSPVS**SP**NNVTLRSSVS**SP**ANINN

Definition repeats

Sequence, long, imperfect, tandem

MRAVVK**SP**IM CHE

KSPSVC**SP**LN

MTSSVC**SP**AG INSVSSTTASF

GSFPVH**SP**IT Q

GTLTC**SP**NV EN

RGSRSH**SP**AH ASN

VGSPLS**SP**L S

MKSSIS**SP**PS HCS

VKSPV**SP**NN VT

LRSSVS**SP**AN INN

Definition repeats

Sequence, long, imperfect, tandem

MRAV**VKSPIM** CHE

KSPSVC**SPLN**

MT**SSVCSPAG** INSVSSTTASF

GSFP**VHSPIT** Q

GTLTC**SPNV** EN

RG**SRSRSHSPA**H ASN

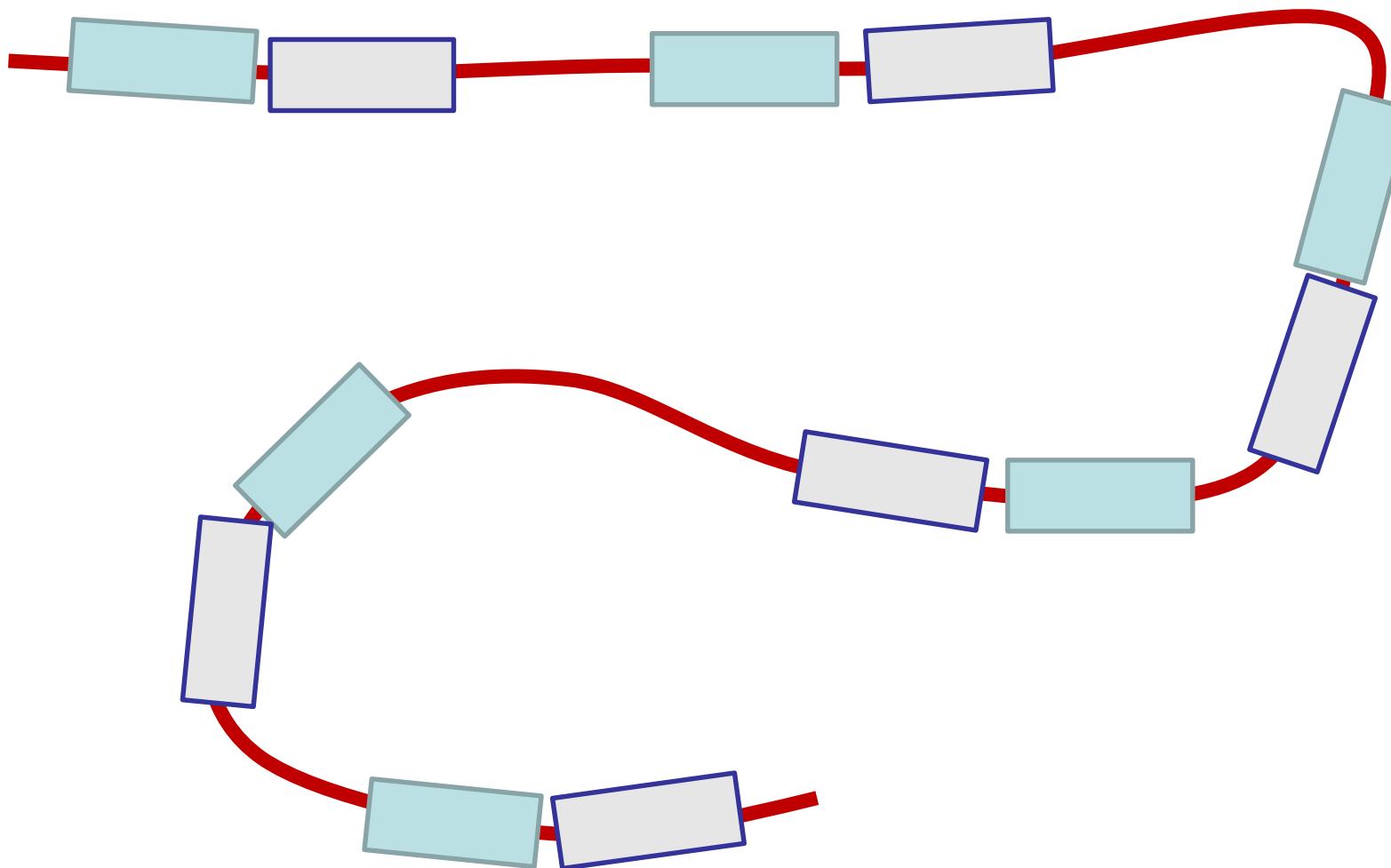
VG**SPLSSP**LS S

MK**SSISSP**PS HCS

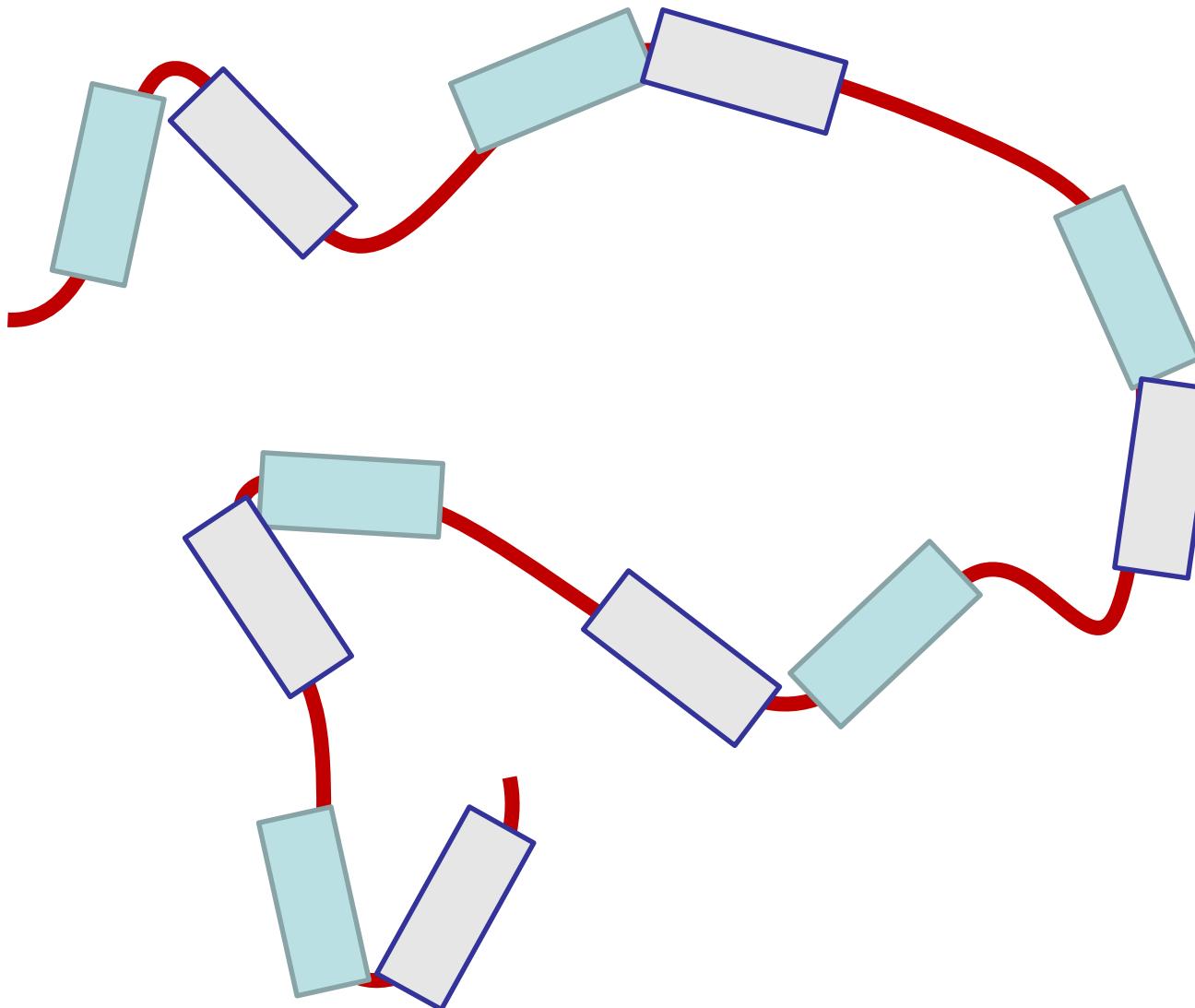
VK**SPVSSP**NN VT

LR**SSVSSP**AN INN

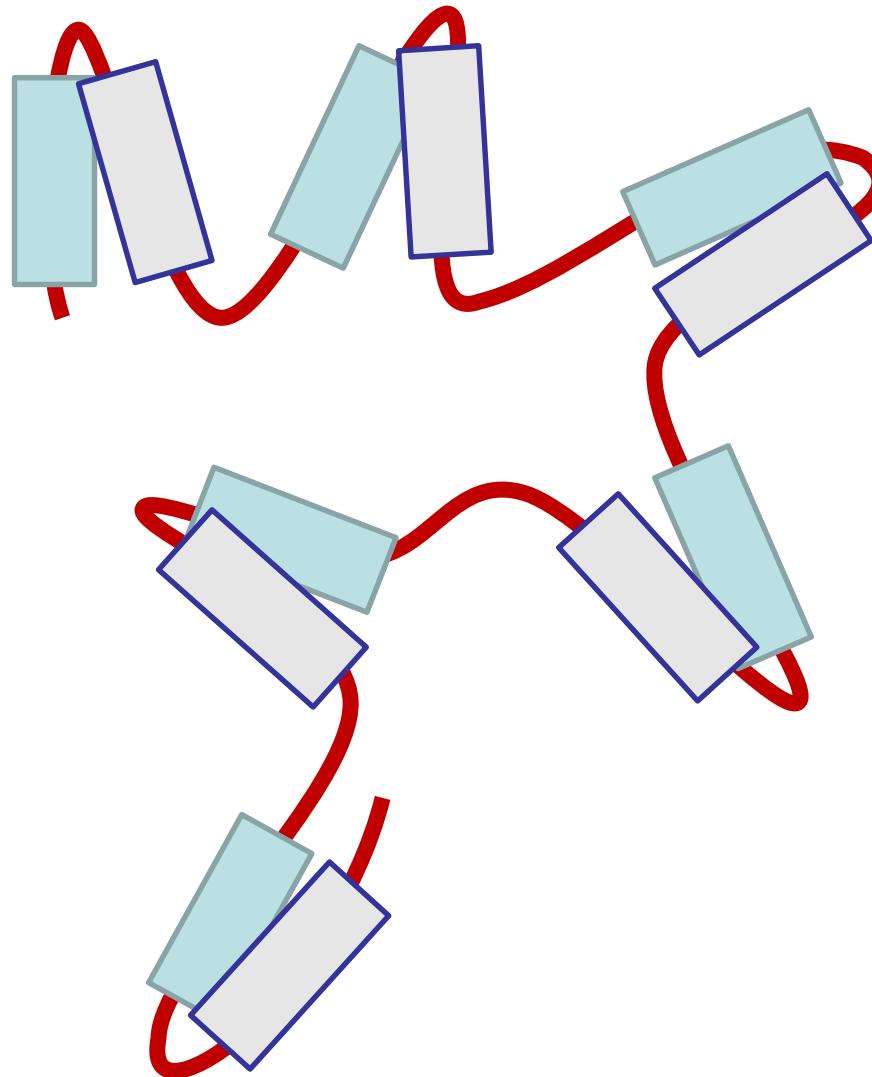
Tandem repeats fold together



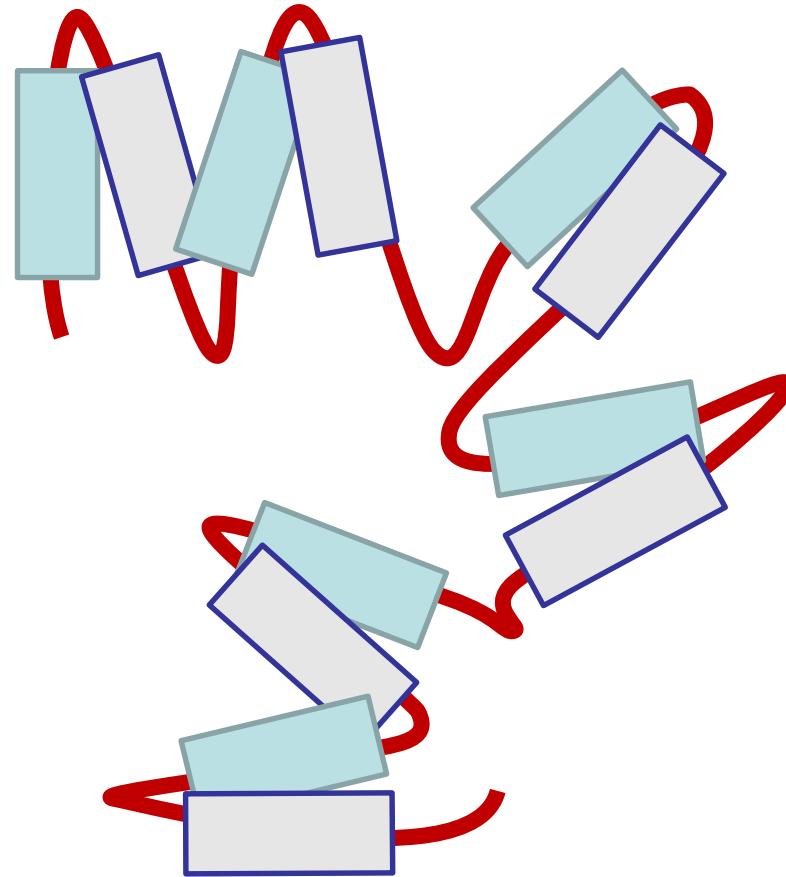
Tandem repeats fold together



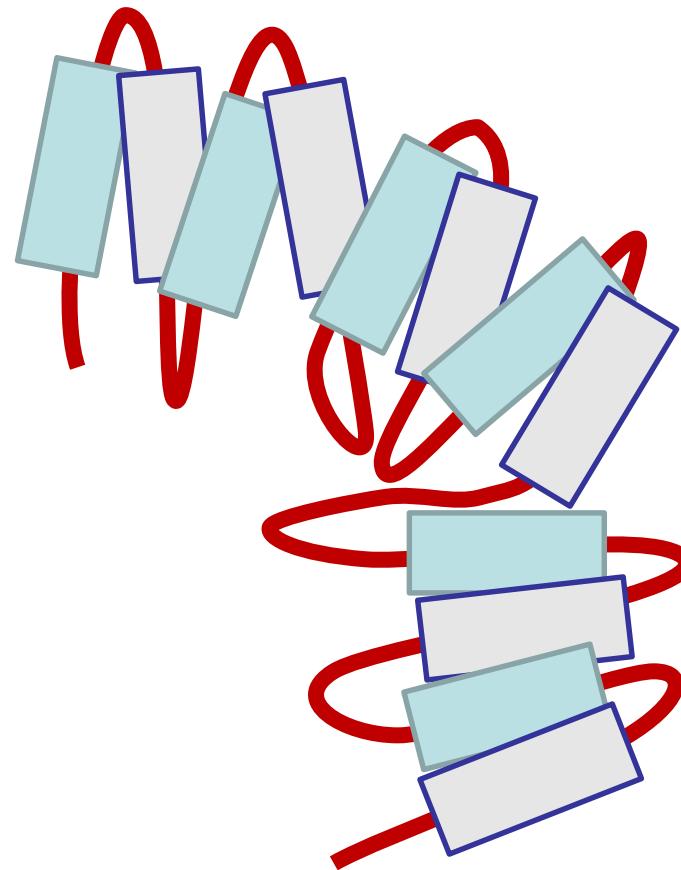
Tandem repeats fold together



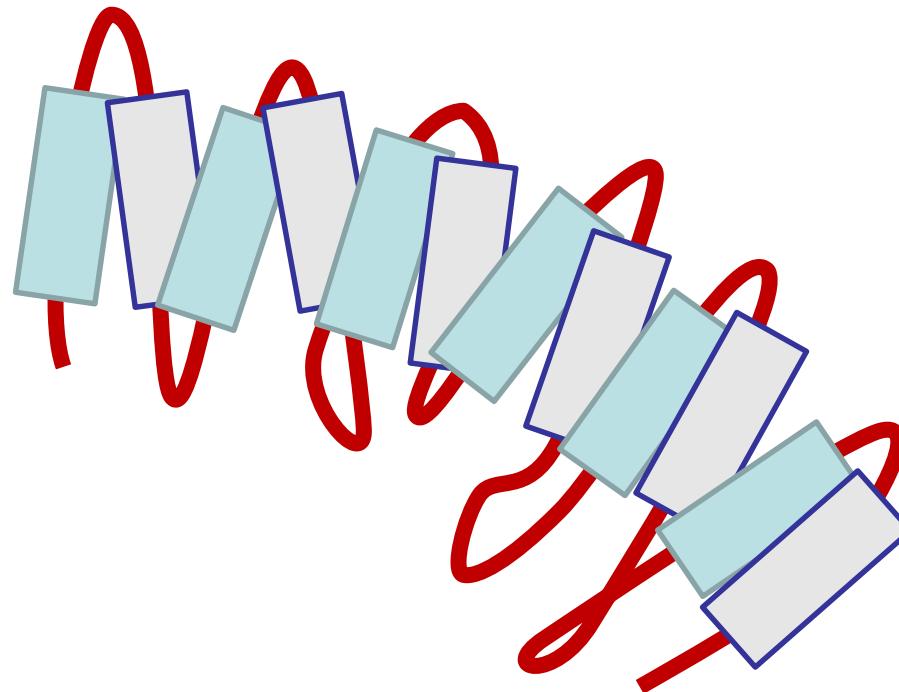
Tandem repeats fold together



Tandem repeats fold together



Tandem repeats fold together



Definition repeats

Sequence, long, imperfect, tandem

MRAV**VKSPIM** CHE

KSPSVC**SPLN**

MT**SSVCSPAG** INSVSSTTASF

GSFP**VHSPIT** Q

GTLTC**SPNV** EN

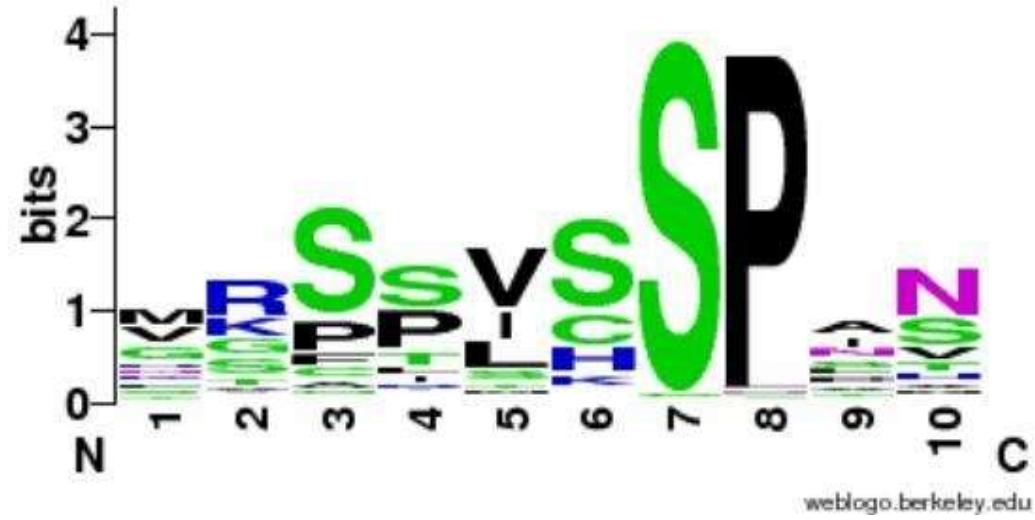
RG**SRSRSHSPA**H ASN

VG**SPLSSP**LS S

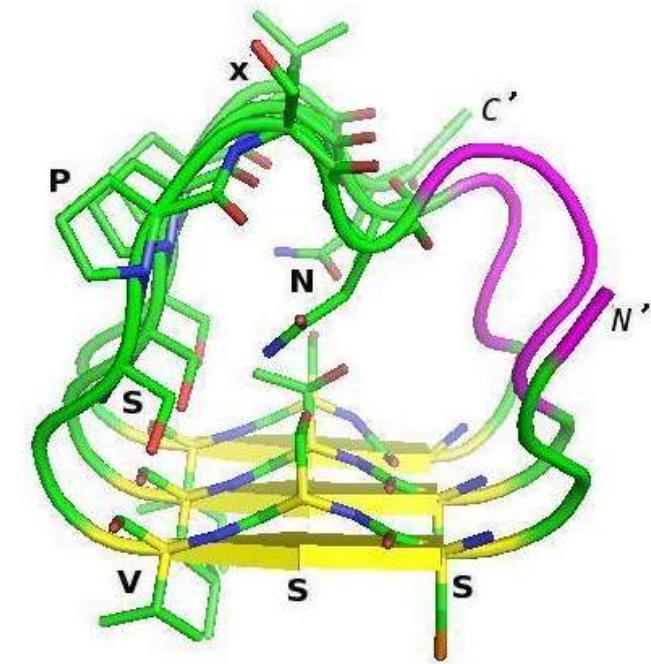
MK**SSISSP**PS HCS

VK**SPVSSP**NN VT

LR**SSVSSP**AN INN



<http://weblogo.berkeley.edu>



(Vlassi et al, 2013)

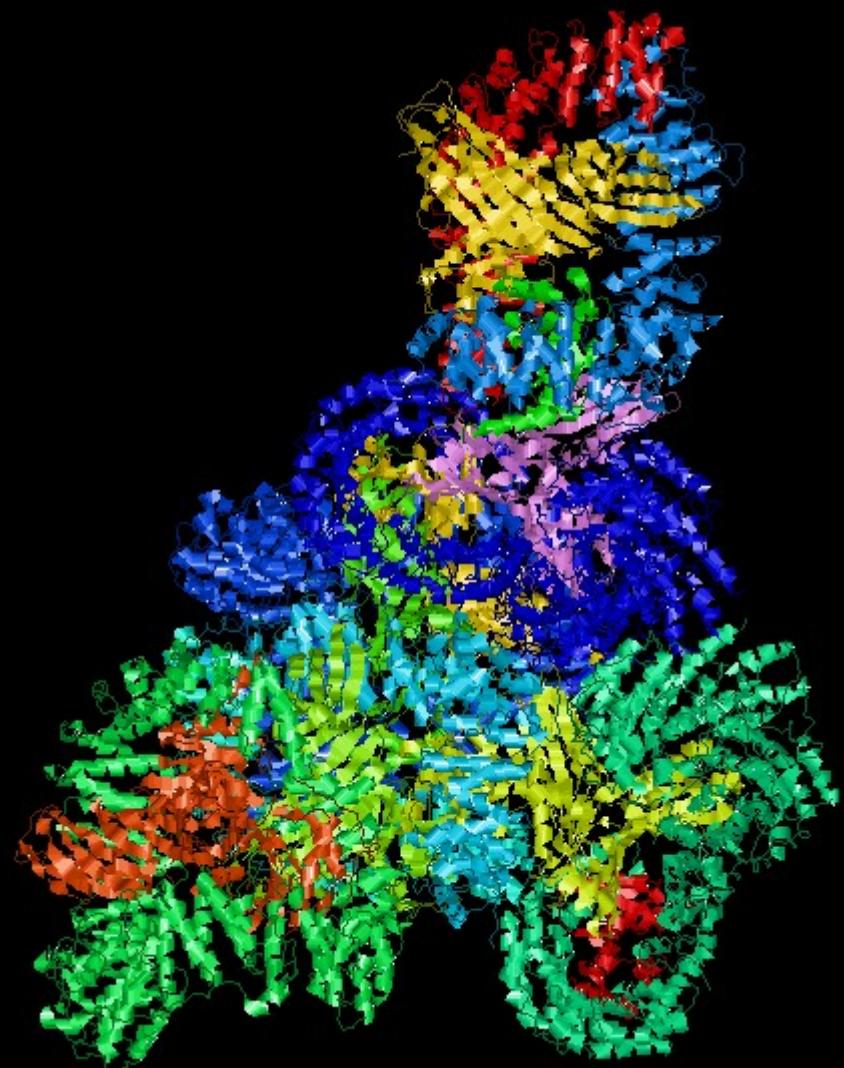
A subunit PP2A structure



PDB:1b3u

Groves et al. (1999) *Cell*

Ap1 Clathrin Adaptor Core

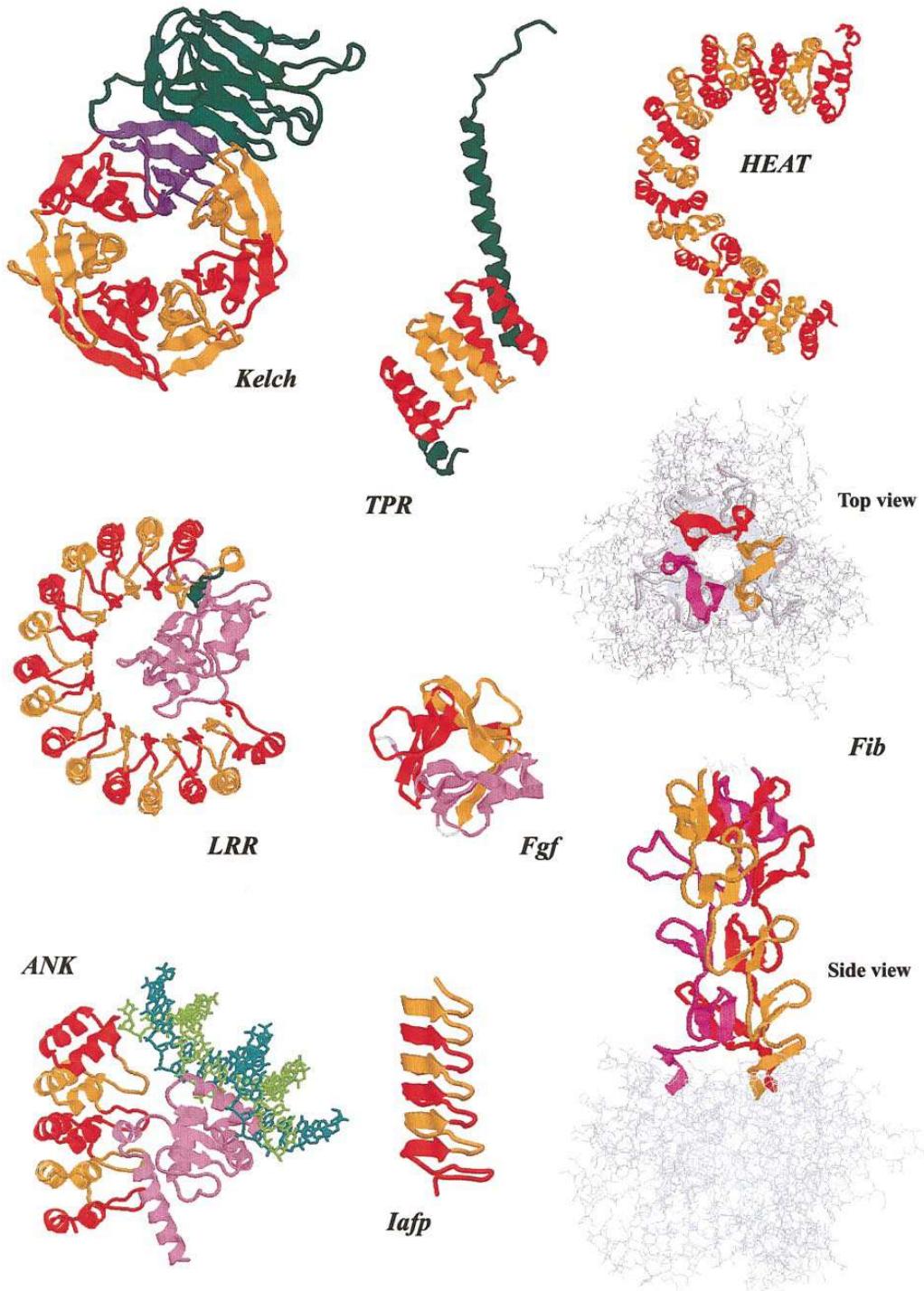


PDB:1w63
Heldwein et al. (2004) *PNAS*

Ap1 Clathrin Adaptor Core



PDB:1w63
Heldwein et al. (2004) *PNAS*



Andrade et al. (2001)
J Struct Biol

Definition CBRs

Perfect repeat: QQQQQQQQQQQQQQ

Imperfect: QQQQPQQQQQQQ

Amino acid type: DDDDDDEEEDEED

Compositionally biased regions (CBRs)

High frequency of one or two amino acids in a region.

Particular case of low complexity region

Detection CBRs

Sometimes straightforward.
N-terminal human Huntingtin.
How many **CBRs** can you find?

```
>sp|P42858|HD_HUMAN Huntingtin OS=Homo sapiens
MATLEKLMKAFESLKSFQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQLPQPPPQAQP
LLPQPQPPPPPPPPPGPAVAEEPLHRPKKELSATKKDRVNHCLTICENIVAQSVRNSPE
FQKLLGIAMELFLLCSDDAESDVRMVADEClnKVlKAlMDSNLPRLQLELYKEIKNGAP
RSLRAALWRFAELAHLVRPQKCRPYLVNLLPCLTRTSKRPEESVQETLAAVPKIMASFG
NFANDNEIKVllKAFIANLKSSSPTIRRTAAGSAVSICQHSRRTQYFYSWllNVllGllV
PVEDEHSTLLlLGvllTLRYLVPllQQVKDTSLKGSGVTRKEMEVSPSAEQLVQVYEL
TLHHTQHQDHNVVTGALEllQQLFRTPPPELLQTLTAVGIGQLTAKEESGGRSRGSI
VELIAGGGSCSPVLSRKQKGKVllGEEEALEDDSESRSVDVSSSALTASVKDEISGELAA
SSGVSTPGSAGHDITEQPRSQHTLQADSVDLASCDLTSSATDGDEEDILSHSSQVSAV
PSDPAMDlNDGTQASSPISDSSQTTEGPDSAVTPSDSSEIVLDGTDNQYLGLQIGQPQD
EDEEATGILPDEASEAFRNSSMALQQAHLLKNMSHCRQPSDSSVDFVLRDEATEPGDQE
NKPCRIKGDIGQSTDSDAPLVHCVRLLSASFLLTGGKNVLVPDRDVRVSVKALALSCVG
AAVALHPESFFSKLYKVPLDTTEYPEEQYVSDILNYIDHGDPQVRGATAILCGTLICSl
```

Detection CBRs

Sometimes straightforward.
N-terminal human Huntingtin.
How many **CBRs** can you find?

```
>sp|P42858|HD_HUMAN Huntingtin OS=Homo sapiens
MATLEKLMKAFESLKSFQQQQQQQQQQQQQQQQQQQQQQPPPPPPPPQLPQPPPQAQP
LLPQPQPPPPPPPPGPAVAEEPLHRPKKELSATKKDRVNHCLTICENIVAQSVRNSPE
FQKLLGIAMELFLLCSDDAESDVRMVADECLNKVIKALMDSNLPRLQLELYKEIKNGAP
RSLRAALWRFAELAHLVRPQKCRPYLVNLLPCLTRTSKRPEESVQETLAAAVPKIMASFG
NFANDNEIKVLLKAFIANLKSSSPTIRRATAAGSAVSICQHSRRTQYFYSWLLNVLLGLLV
PVEDEHSTLLILGVLLTRYLVPOLLQQVKDTSLKGSGVTRKEMEVSPSAEQLVQVYEL
TLHHTQHQDHNVVTGALELLQQLFRTPPPELLQTLTAVGGIGQLTAKEESGGRSRSGSI
VELIAGGGSCSPVLSRKQKGKVLLGEEEALEDDSESRSVDVSSSALTASVKDEISGELAA
SSGVSTPGSAGHDITEQPRSQHTLQADSDVLASCDLTSSATDGDEEDILSHSSSQVSAV
PSDPAMDLDGTQASSPISDSSQTTEGPDSAVTPSDSSEIVLDGTDNQYLGLQIGQPQD
EDEEATGILPDEASEAFRNSSMALQQAHLLKNMSHCRQPSDSSVDFVLRDEATEPGDQE
NKPCRIKGDIGQSTDSSAPLVHCVRLLSASFLLTGGKNVLVPDRDVRVSVKALALSCVG
AAVALHPESFFSKLYKVPLDTTEYPEEQYVSDILNYIDHGDPQVRGATAILCGTLICSL
```

Detection CBRs

Sometimes straightforward.
N-terminal human Huntingtin.
How many **CBRs** can you find?

```
>sp|P42858|HD_HUMAN Huntingtin OS=Homo sapiens
MATLEKLMKAFESLKSFQQQQQQQQQQQQQQQQQQQQQQPPPPPPPQQLPQPPPQAQP
LLPQPQPPPPPPPPGPAVAEELHRPKKELSATKKDRVNHCLTICENIVAQSVRNSPE
FQKLLGIAMELFLLCSDDAESDVRMVADECINKVIKALMDSNLPRLQLEYKEIKNGAP
RSLRAALWRFAELAHLVRPQKCRPYLVNLLPCLTRTSKRPEESVQETLAAVPKIMASFG
NFANDNEIKVLLKAFIANLKSSSPTIRRATAAGSAVSICQHSRRTQYFYSWLLNVLLGLLV
PVEDEHSTLLILGVLLTLRYLVPPLLQQVKDTSLKGSGFGVTRKEMEVSPSAEQLVQVYEL
TLHHTQHQDHNVVTGALELLQQLFRTPPPELLQTLTAVGGIGQLTAKEESGGRSRSGSI
VELIAGGGSCSPVLSRKQKGKVLLGEEEAEDDSERSRDVSSSALTASVKDEISGELAA
SSGVSTPGSAGHDITEQPRSQHTLQADSDLASCDLTSSATDGDEEDILSHSSQVSAV
PSDPAMDLDGTQASSPISDSSQTTEGPDSAVTPSDSSEIVLDGTDNQYLGLQIGQPQD
EDEEATGILPDEASEAFRNSSMALQQAHLLKNMSHCRQPSDSSVDKFVLRDEATEPGDQE
NKPCRIKGDIGQSTDSSAPLVHCVRLLSASFLLTGGKNVLVPDRDVRVSVKALALSCVG
AAVALHPESFFSKLYKVPLDTTEYPEEQYVSDILNYIDHGDPQVRGATAILCGTLICSL
```

Detection CBRs

Sometimes straightforward.
N-terminal human Huntingtin.
How many **CBRs** can you find?

>sp|P42858|HD_HUMAN Huntingtin OS=Homo sapiens
MATLEKLMKAFESLKSFQQQQQQQQQQQQQQQQQQQQQQQQQQPPPPPPPPPPPQLPQPPPQAQP
LLPQPQPPPPPPPPGPAVAEELHRPKKELSATKKDRVNHCLTICENIVAQSVRNSPE
FQKLLGIAMELFLLCSDDAEDVCRMVADECLNKVIKALMDSNLPRQLELYKEIKKNGAP
RSLRAALWRFAELAHLVRPQKCRPYLVNLLPCLTRTSKRPEESVQETLAAAVPKIMASFG
NFANDNEIKVLLKAFIANLKSSSPTIRRTAAGSAVSICQHSRRTQYFYSWLLNVLLGLLV
PVEDEHSTLLILGVLLTLRYLVPLLQQQVKDTSLKGSFGVTRKEMEVSPSAEQLVQVYEL
TLHHTQHQDHNVVTGALELLQQLFRTPPPELLQTLTAVGGIGQLTAKEESGGRSRSGSI
VELIAGGGSSCSPVLSRKQKGKVLLGEEEALEDDESESRSDVSSALTASVKDEISGELAA
SSGVSTPGSAGHDIITEQPRSQHTLQADSVDLASCDLTSSATDGDEEDILSHSSSQVSAV
PSDPAMDLNQASSPISDSSQTTTEGPDSAVENTPSDSSEIVLDGTDNQYLGLQIGQPQD
EDEEATGILPDEASEAFRNSSMALQQAHLLKNMSHCRQPSDSSVDKFVLRDEATEPGDQE
NKPCRIKGDIGQSTDSSAPLVHCVRLLSASFLLTGGKNVLVPDRDVRVSVKALALSCVG
AAVALHPESFFSKLYKVPLDTTEYPEEQYVSDILNYIDHGDPQVRGATAILCGTLICSI

Detection repeats

Sometimes straightforward.
N-terminal human Huntingtin.
How many **repeats** can you find?

```
>sp|P42858|HD_HUMAN Huntingtin OS=Homo sapiens
MATLEKLMKAFESLKSFQQQQQQQQQQQQQQQQQQQQQQQQQQQQQLPQPPPQAQP
LLPQPQPPPPPPPPPGPAVAEEPLHRPKKELSATKKDRVNHCLTICENIVAQSVRNSPE
FQKLLGIAMELFLLCSDDAESDVRMVADEClnKVlKAlMDSNLPRLQLELYKEIKNGAP
RSLRAALWRFAELAHLVRPQKCRPYLVNLLPCLTRTSKRPEESVQETLAAVPKIMASFG
NFANDNEIKVllKAFIANLKSSSPTIRRTAAGSAVSICQHSRRTQYFYSWllNVllGllV
PVEDEHSTLLlLGvllTLRYLVPllQQVKDTSLKGSGVTRKEMEVSPSAEQLVQVYEL
TLHHTQHQDHNVVTGALEllQQLFRTPPPELLQTLTAVGIGQLTAKEESGGRSRGSI
VELIAGGGSCSPVLSRKQKGKVllGEEEALEDDSESRSVDVSSSALTASVKDEISGELAA
SSGVSTPGSAGHDlITEQPRSQHTLQADSVDLASCDLTSSATDGDEEDILSHSSQVSAV
PSDPAMDlNDGTQASSPISDSSQTTEGPDSAVTPSDSSEIVLDGTDNQYLGLQIGQPQD
EDEEATGILPDEASEAFRNSSMALQQAHLLKNMSHCRQPSDSSVDFVLRDEATEPGDQE
NKPCRIKGDIGQSTDSDAPLVHCVRLLSASFLLTGGKNVLVPDRDVRVSVKALALSCVG
AAVALHPESFFSKLYKVPLDTTEYPEEQYVSDILNYIDHGDPQVRGATAILCGTLICSl
```

Detection repeats

Often NOT straightforward.

N-terminal human Huntingtin.

How many **repeats** can you find?

```
>sp|P42858|HD_HUMAN Huntingtin OS=Homo sapiens
MATLEKLMKAFESLKSFQQQQQQQQQQQQQQQQQQQQQQQQQQQQQLPQPPPQAQP
LLPQPQPPPPPPPPPGPAVAEEPLHRPKKELSATKKDRVNHCLTICENIVAQSVRNSPE
FQKLLGIAMELFLLCSDDAESDVRMVADEClnkVIKALMDSNLPRLQLEYKEIKNGAP
RSLRAALWRFAELAHLVRPQKCRPYLVNLLPCLTRTSKRPEESVQETLAAAVPKIMASFG
NFANDNEIKVLLKAFIANLKSSSPTIRRTAACGSAVSICQHSRRTQYFYSWLLNVLLGLLV
PVEDEHSTLLILGVLLTRYLVPLLQQVKDTSLKGSGVTRKEMEVSPPSAEQLVQVYEL
TLHHTQHQDHNVVTGALELLQQLFRTPPPELLQTLTAVGGIGQLTAKEESGGRSRSGSI
VELIAGGGSCSPVLSRKQKGKVLLGEEEALEDDSESRSVDVSSSALTASVKDEISGELAA
SSGVSTPGSAGHDITEQPRSQHTLQADSVLASCDLTSSATDGDEEDILSHSSQVSAV
PSDPAMDLDGTQASSPISDSSQTTEGPDSAVTPSDSSEIVLDGTDNQYLGLQIGQPQD
EDEEATGILPDEASEAFRNSSMALQQAHLLKNMSHCRQPSDSSVDFVLRDEATEPGDQE
NKPCRIKGDIGQSTDSSAPLVHCVRLLSASFLLTGGKNVLVPDRDVRVSVKALALSCVG
AAVALHPESFFSKLYKVPLDTTEYPEEQYVSDILNYIDHGDPQVRGATAILCGTLICSL
```

Detection repeats

Often NOT straightforward.

N-terminal human Huntingtin.

How many **repeats** can you find?

EFQKLLGIAMELFLLCSD**DA**ESDVRMVADECLNKVIKA
CRPYLVNLLPCLRTSKR**P-EESVQETLAAAVPKIMAS**
NDNEIKVLLKAFIANLKS**SSPTIRRTAAGSAVSICQHS**
TQYFYSWLLNVLLGLLVP**VEDEHSTLLILGVLLTLRYL**
PSAEQLVQVYELTLHHTQ**HQDHNVVTGALELLQQLFRT**

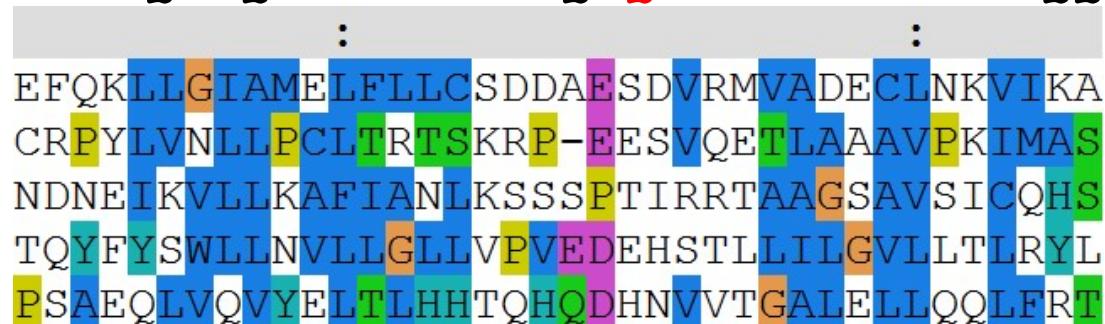
Detection repeats

Often NOT straightforward.

N-terminal human Huntingtin.

How many **repeats** can you find?

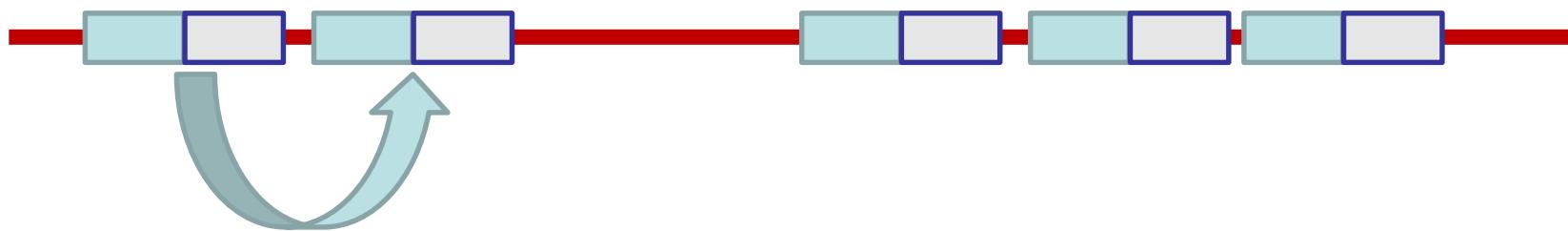
EFQKLLGIAMELFLLCSDDAESDVRMVADECLNKVIKA
CRPYLVNLLPCLTRTSKR~~P~~-EESVQETLAAAVPKIMAS
NDNEIKVLLKAFIANLKS~~SSPTIRRTAAGSAVSICQHS~~
TQYFYSWLLNVLLGLLVP~~VEDEHSTLLILGVLLTLRYL~~
PSAEQLVQVYELTLHHTQ~~HQDHNVVTGALELLQQLFRT~~



Detection of repeats

Dotplots

Comparing a sequence against itself



Detection of repeats

Dotplots

TLRSSVSSSPANINNS
NMTSSVCSPANISV

Detection of repeats

Dotplots



TLRSSVSSSPANINNS
|
NMTSSVCSPANISV

1 match

Detection of repeats

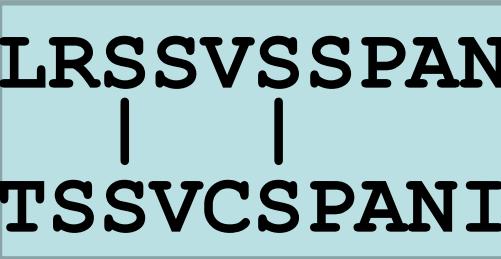
Dotplots

TLRSSVSSSPANINNS
| | | | | | | |
NMTSSVCSPANISV

8 matches

Detection of repeats

Dotplots



TLRSSVSSSPANINNS
| |
NMTSSVCSPANISV

2 matches

Detection of repeats

Dotplots

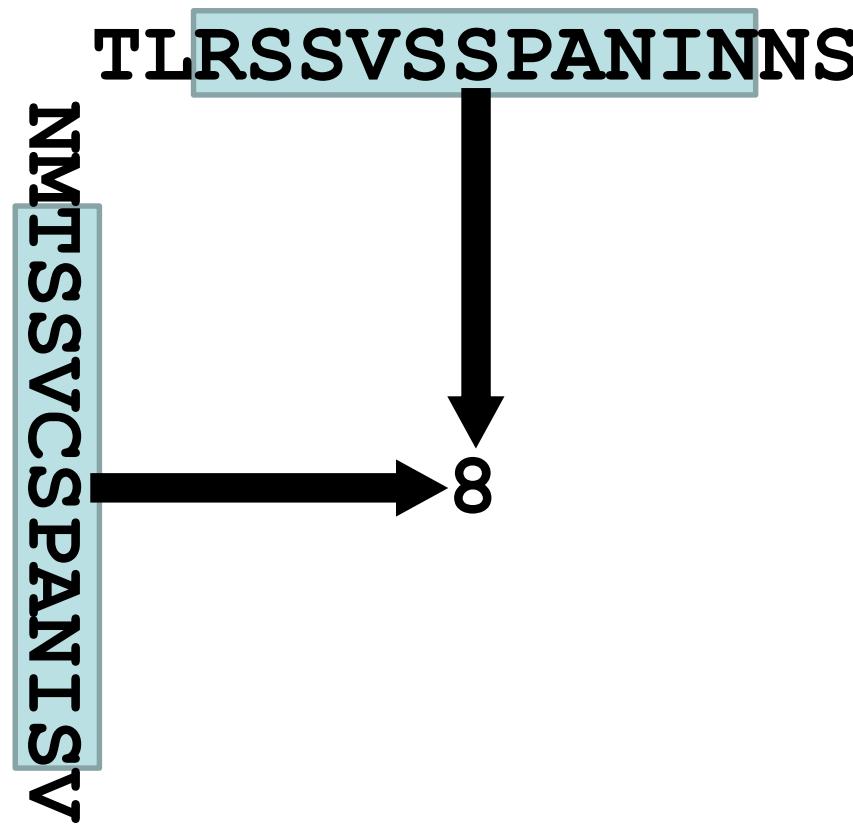


TLRSSVSSSPANINNS
|
NMTSSVCSPLANISV

1 match

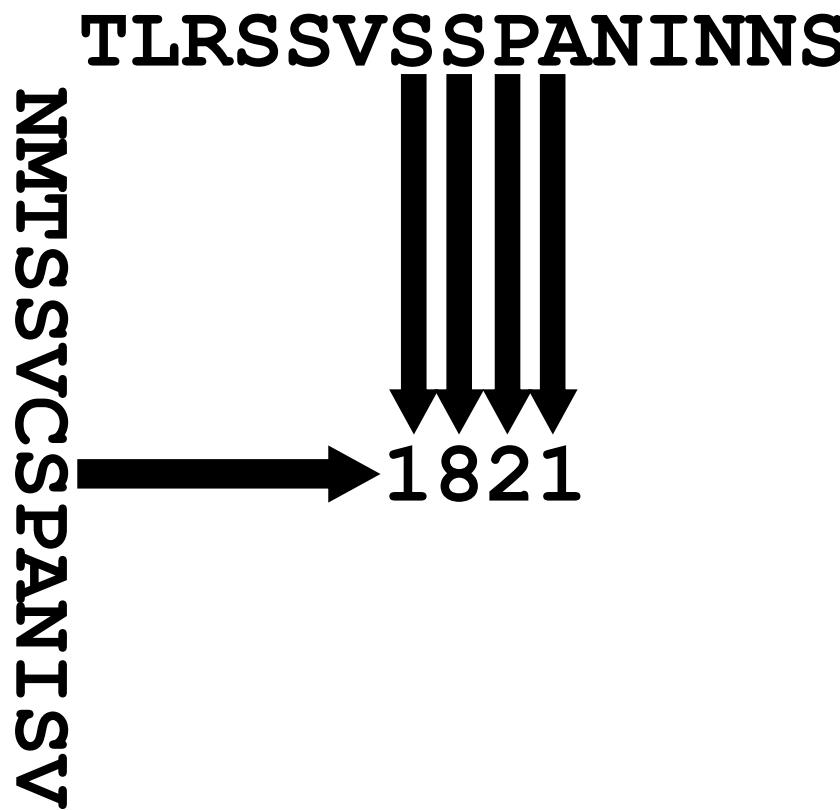
Detection of repeats

Dotplots



Detection of repeats

Dotplots



SEQUENCE 1

SEQUENCE 2

Window size

15

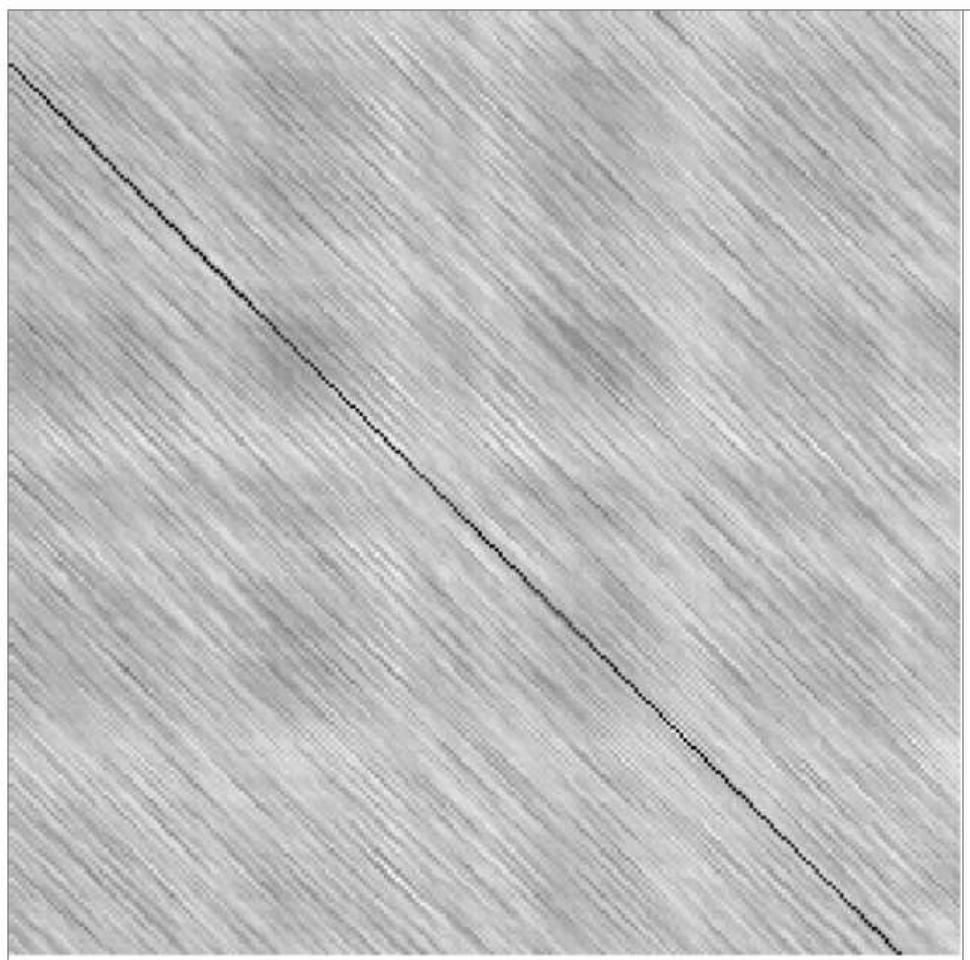
Scoring matrix

BLOSUM 62

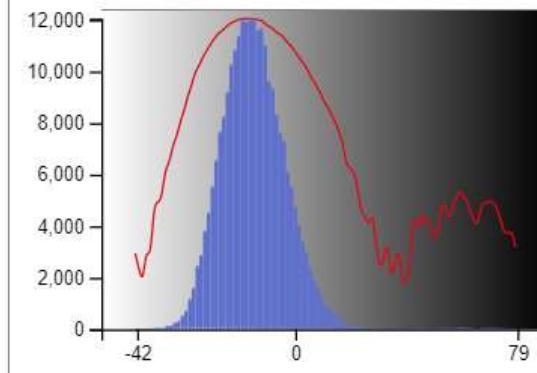
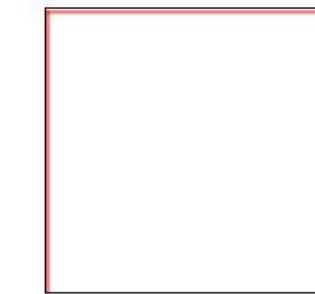


Sequence 1

Sequence 2



[246 x 244] # Score at (1:M, 1:R) : -8



Seq1:1

[MTMDKSEL]VQKAKLAEQAERYDDM**A**AMKAVTEQGHE
RKPLQTPT**P**IRRLWTMDTSELVQ**K**AKLAEQAERYDDM

Exercise 1. Using Dotlet with the human mineralocorticoid receptor (MR)

- Go to the Dotlet web page:

<http://dotlet.vital-it.ch>

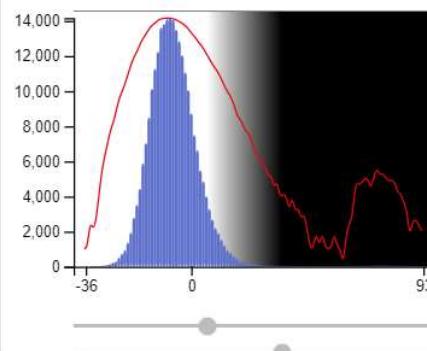
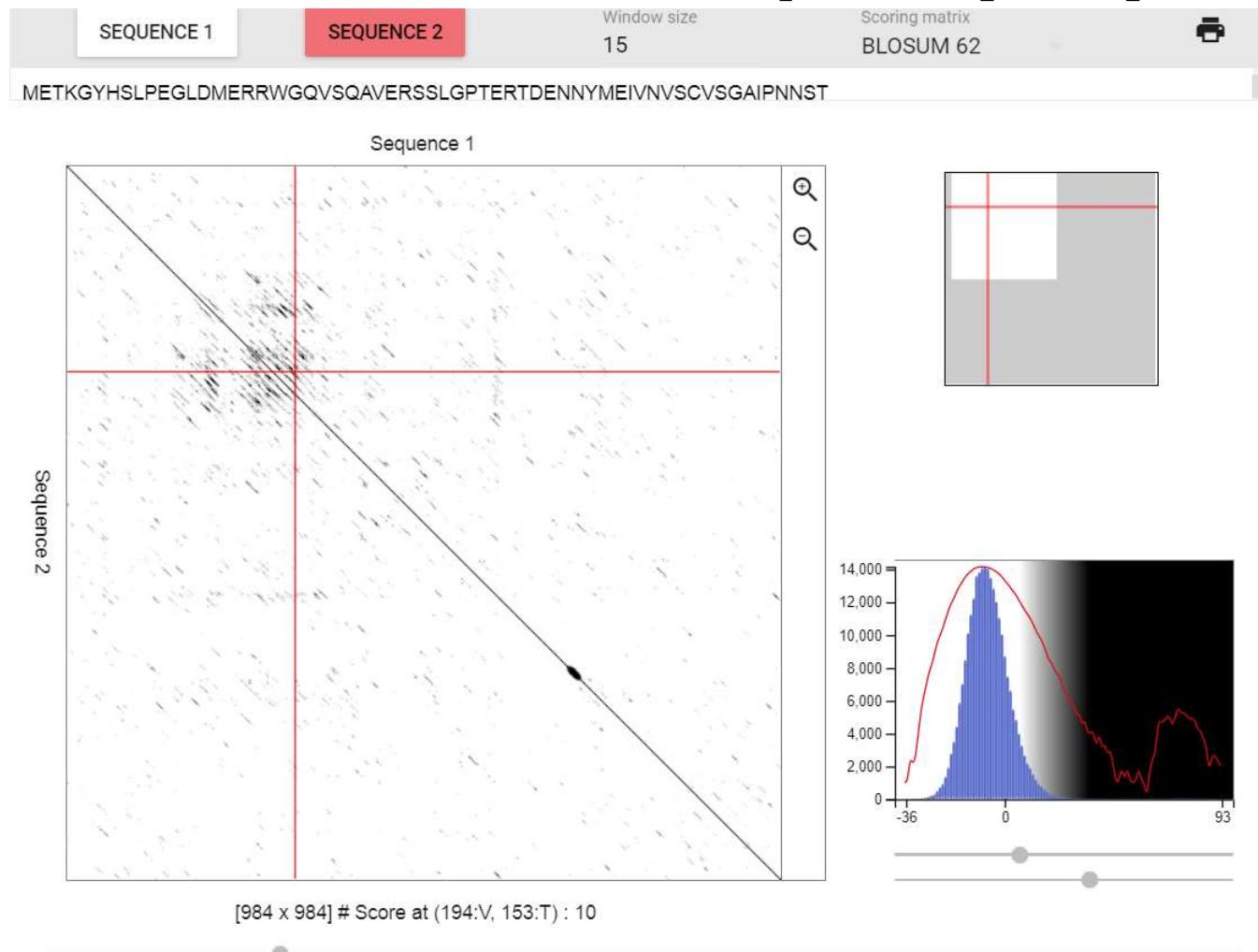
- Click on the input button and paste the sequence of the human mineralocorticoid receptor (UniProt id P08235)

- Click on the “compute” button

- Try to find combinations of parameters that show patterns in the dot plot (Hint: You can adjust this finely using the arrows)

- Find repetitions clicking in the diagonal patterns

Exercise 1. Using Dotlet with the human mineralocorticoid receptor (MR)



Seq1:194

NTPLRSFMSD**S**G**S**VNNGGVRAVVKSPIM**C**HEKSP**S**V**C**SPL**N**T**I**SSVCSPAGIN**S**V**S**TTASFGSFPVHSPIT
YSYEQQNQQG**S****M**PAKIYQNVEQLVKFYK**G****N**GHRP**S**T**L****S**CV**N**TPLRSFMSD**S**G**S**VNNGGVRAVVKSPIM**C**HE

Detection of repeats

Using a multiple sequence alignment helps.
Conserved repeated patterns

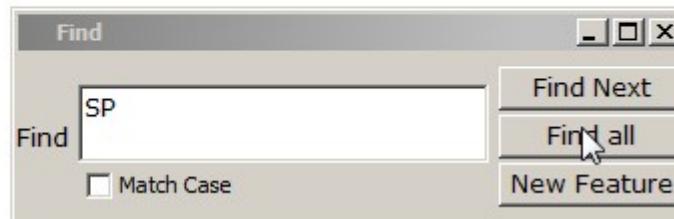
	240	250	260	270	280	290	300	310	320
<i>mr_danio/1-970</i>	---	TYF	-DSDCP	TLD SAT SSLTHCQHTSPN I CSPVKSSIVGS PPLPSPLSVMKSPVSSPHSIGSVRSPLSC	-	-NTNMRSSVSSPTTNG			
<i>mr_nutilus/1-973</i>	---	TYF	-DSDCP	-SLASASTNL TQGHHTSPNTCS PVKSSMVGS PPLSVMKSPVSSPRSIGSVRSPLSC	-	-NTNMRSSVSSPTTNG			
<i>mr_cyprinus/1-971</i>	---	TFF	-DSDCP	-SLASHTHNL IQGQHTSPNTCS PVKSSVVGS PPLSVMKSPVSSPHSIGSVSSPLSC	-	-NTNMRSSVSSPTTYG			
<i>mr_oryzias/1-994</i>	TCFGPQCSAVSSPVQSQTSCAATLANI	KRRNSVTCS PVESCTVGS PPLTSPLN	IMRSPMSSPHSMSSVRSPSCSTTCNI RSSVSSPT	---					
<i>mr_takifugu/1-991</i>	MCF GPMCSSVSSPVQSQTSCA STLPN I	KRRNSATCSPV ESS TVGS PPLTSPLN	IMRSP I SSPQSMSSVRSPSCSTTSNI RSSVSSPT	---					
<i>mr_oreochromis/1-994</i>	TCFAPLCSSVSSPVQSQTSCAATLANI	KRRNSVTCS PVESSTVGS PPLTSPLNVMRSPMSSPQSMSSVRSPSCSTTCNI RSSVSSPT	---						
<i>mr_xenopus/1-979</i>	-- FGNF	-- TVHSPVNQVTPKSCS PHTDNRC SIAHSP	-- AGTVES	PLSSPVSSMRSPISSPPSHASL KSPVSSPNNITVRPSVSSPGNI					
<i>mr_anolis/1-990</i>	-- FGNF	-- TVSSPVNQGTPLSCSPNIE NRGSMLHSPPHASNMGS	PLSSP I	SSMKSPISSPPSHCSVKSPVSSPNNITMRSSVSSPANM					
<i>mr_alligator/1-985</i>	-- FGNF	-- VVNSPINQGTPLSCSPNIE NRGSMLHS PAHASNVGS	PLSSP I	SSMKSPISSPPSHCSVKSPVSSPNNITMRSSVSSPANM					
<i>mr_taeniopygia/1-981</i>	-- FGNF	-- SMHSPMGQGTPLSRSPN VENRG SMLHS PAHIS NVGS	PLSSP I	SSMKSPISSPPSHCSVKSPVSSPNNITMRSSVSSPANL					
<i>mr_gallus/1-986</i>	-- FGNF	-- AMHSP I GQGTPLSRSPN VESRG SMLHS PAHV S NVGS	PLSSP I	SSMKSPISSPPSHCSVKSPVSSPNNITMRSSVSSPANM					
<i>mr_monodelphis/1-993</i>	-- FG SF	-- PVHSP ITQGTPLPCSPN VENRG S RSHSPV HASNVGS	PLSSP I	SSMKSPISSPPSHCSVKSPVSSPNNVTMRSSVSSPANIN					
<i>mr_mus/1-980</i>	-- FG SF	-- PVHSP ITQGTSLTCSPS VENRG S RSHSPV HASNVGS	PLSSP LSSM	SKSPISSPPSHCSVKSPVSSPNNVPLRSSVSSPANLN					
<i>mr_rattus/1-981</i>	-- FG SF	-- PVHSP ITQGTSLTCSPS VENRG S RSHSPV HASNVGS	PLSSP LSSM	KSPISSPPSHCSVKSPVSSPNNVPLRSSVSSPANLN					
<i>mr_homo/1-984</i>	-- FG SF	-- PVHSP ITQGTPLTCS PNAENRG S RSHSPV HASNVGS	PLSSP LSSM	KSSISSPPSHCSVKSPVSSPNNVTLRSSVSSPANIN					
<i>mr_equus/1-984</i>	-- FGNF	-- TVHSP ITQGTPLCS PNVENRG S RSHSPV HASNVGS	PLSSP LSSM	KSPISSPPSHCSVKSPVSSPNNVTLRSSVSSPANIN					

JalView with Regular Expression searches

Detection of repeats

Using a multiple sequence alignment helps
Conserved repeated patterns

	240	250	260	270	280	290	300	310	320	
<i>mr_danio/1-970</i>	---	TYF	-	DSDCP	TLD SAT SSLTHCQHTSPN	I C SPVKSSIVGSPPL	P SPL SVMKSPVSS	S PHSIGSVRSPLSC	-	NTNM RSS VSSP TTNG
<i>mr_utilis/1-973</i>	---	TYF	-	DSDCP	- SLASASTNL TQGHHTSPN	TCS PVKSSMVGS	PPL ASPL SVMKSPVSS	SPRSIGSVRSPLSC	-	NTNM RSS VSSP TTNG
<i>mr_cyprinus/1-971</i>	---	TFF	-	DSDCP	- SLASATHTNL IQGQHTSPN	TCS PVKSSVVGS	PPL ASPL SVMKSPVSS	SPHSIGSVSSPLSC	-	NTNM RSS VSSP TTNG
<i>mr_oryzias/1-994</i>	TCFGPQCSAVSSPVSQTSCAATLANI	KRRNSVT	KRRNSATCSPV	EESSTVGSPPLTSP	LNIMRSPMSSPHSMSSVR	SPSCSTTCNI RSSVSSP	T	---	---	---
<i>mr_takifugu/1-991</i>	MCFGP MCSSVSSPVSQTSCAATLPI	KRRNSATCSPV	EESSTVGSPPLTSP	LNIMRSPMSSPHSMSSVR	SPSCSTTSNI RSSVSSP	T	---	---	---	---
<i>mr_oreochromis/1-994</i>	TCFA PLCSSVSSPVSQTSCAATLANI	KRRNSVT	KRRNSATCSPV	EESSTVGSPPLTSP	LNVMRSPMSSPQSMSSVR	SPSCSTTCNI RSSVSSP	T	---	---	---
<i>mr_xenopus/1-979</i>	--- FGNF -	TVHSPVNQVTPKSCSP	H TDNRC	SPHTDNRC	IAHSP -	- AGTVES	PLSSPVSSMRSP	I SPPSHASL	KSPVSSP	NNITVRPSVSSP
<i>mr_anolis/1-990</i>	--- FGNF -	TVSSPVNQGTP	LSCSPN	I ENRG	SMLHSPPHASNM	GSP	PLSSPIISSMKSP	I SPPSHCSVK	SPVSSP	NNITMRSSVSSPANM
<i>mr_alligator/1-985</i>	--- FGNF -	VVN	SPINQGTP	LSCSPN	I ENRG	SMLHSPPHASNV	GSP	PLSSPIISSMKSP	I SPPSHCSVK	SPVSSP
<i>mr_taeniopygia/1-981</i>	--- FGNF -	SMHSPM	GQGTPLSRSP	NVENRG	SMLHSPPAHISNV	GSP	PLSSPIISSMKSP	I SPPSHCSVK	SPVSSP	NNITMRSSVSSPANL
<i>mr_gallus/1-986</i>	--- FGNF -	AMHSP	I GQGTPLSRSP	NVESRG	SMLHSPPAHISNV	GSP	PLSSPIISSMKSP	I SPPSHCSVK	SPVSSP	NNITMRSSVSSPANM
<i>mr_monodelphis/1-993</i>	--- FG SF -	PVHSP	ITQGTPLPCSP	NVENR	SSVSHSPPPAHISNV	GSP	PLSSPIISSMKSP	I SPPSHCSVK	SPVSSP	NNVTMRSSVSSPANIN
<i>mr_mus/1-980</i>	--- FG SF -	PVHSP	ITQGTSLTCSP	SVENRG	SRSHSPPPAHISNV	GSP	PLSSPLSSMKSP	I SPPSHCSVK	SPVSSP	NNVPLRSSVSSPANLN
<i>mr_rattus/1-981</i>	--- FG SF -	PVHSP	ITQGTSLTCSP	SVENRG	SRSHSPPPAHISNV	GSP	PLSSPLSSMKSP	I SPPSHCSVK	SPVSSP	NNVPLRSSVSSPANLN
<i>mr_homo/1-984</i>	--- FG SF -	PVHSP	ITQGTPLTCSP	NAENRG	SRSHSPPPAHISNV	GSP	PLSSPLSSMKSS	I SPPSHCSVK	SPVSSP	NNVTLRSSVSSPANIN
<i>mr_equus/1-984</i>	--- FGNF -	TVHSP	ITQGTPLTCSP	NVENRG	SRSHSPPPAHISNV	GSP	PLSSPLSSMKSP	I SPPSHCSVK	SPVSSP	NNVTLRSSVSSPANIN



JalView with Regular Expression searches

Detection of repeats

Using a multiple sequence alignment helps
Conserved repeated patterns

JalView with Regular Expression searches

Detection of repeats

Using a multiple sequence alignment helps
Conserved repeated patterns

JalView with Regular Expression searches

- Regular Expressions:

[L|S]P.A

matches L or S, followed by P, followed by anything, followed by A

Detection of repeats

Using a multiple sequence alignment helps
Conserved repeated patterns

JalView with Regular Expression searches

- Regular Expressions:

[L S] P . A

matches L or S, followed by P, followed by anything, followed by A

Which one is not matched?

- LPTA, SPAA, LPAA, LPAP, SPLA

Detection of repeats

Using a multiple sequence alignment helps
Conserved repeated patterns

JalView with Regular Expression searches

- Regular Expressions:

[L S] P . A

matches L or S, followed by P, followed by anything, followed by A

Which one is not matched?

- LPTA, SPAA, LPFA, **LPAP**, SPLA

Exercise 2. Using JalView with a MSA of the MR with orthologs

- Load the multiple sequence alignment of the MR in JalView: MR1_fasta.txt (from URL: https://cbdm.uni-mainz.de/files/2015/02/MR1_fasta.txt)
- Use the “Select > find” (of Ctrl+F) option with a regular expression and mark all matches (**click the “Find all” option!**)
- Try to find the expression that matches more repeats. How many repeats do you see? How long are they? Would you correct the alignment based on these findings?

|158508572|Hsapie
 |31324675|Cjacchus
 |126331313|Mdomestica
 |73978292|Clupus
 |301763180|Amelanoleuca
 |6981208|Rnorvegicus
 |144227212|Mmusculus
 |148224443|Xlaevis
 |327274009|Acarolinensis
 |115529242|Tguttata
 |225936142|Ggallus
 |239923135|Rrutilus
 |154240734|Drerio

#T1	#T2	#T3	#T4	#T5	#T6	#T7
..170... ...180... ...190... ...200... ...210... ...220... ...230... ...240... ...250... ...260...

#F1

#F2

#F3

#F4

|158508572|Hsapie
 |31324675|Cjacchus
 |126331313|Mdomestica
 |73978292|Clupus
 |301763180|Amelanoleuca
 |6981208|Rnorvegicus
 |144227212|Mmusculus
 |148224443|Xlaevis
 |327274009|Acarolinensis
 |115529242|Tguttata
 |225936142|Ggallus
 |239923135|Rrutilus
 |154240734|Drerio

#T8	#T9	#T10	#T11	#T12	#T13	#T14	#T15
..270... ...280... ...290... ...300... ...310... ...320... ...330... ...340... ...350... ...360...

#F5

#F6

#F7

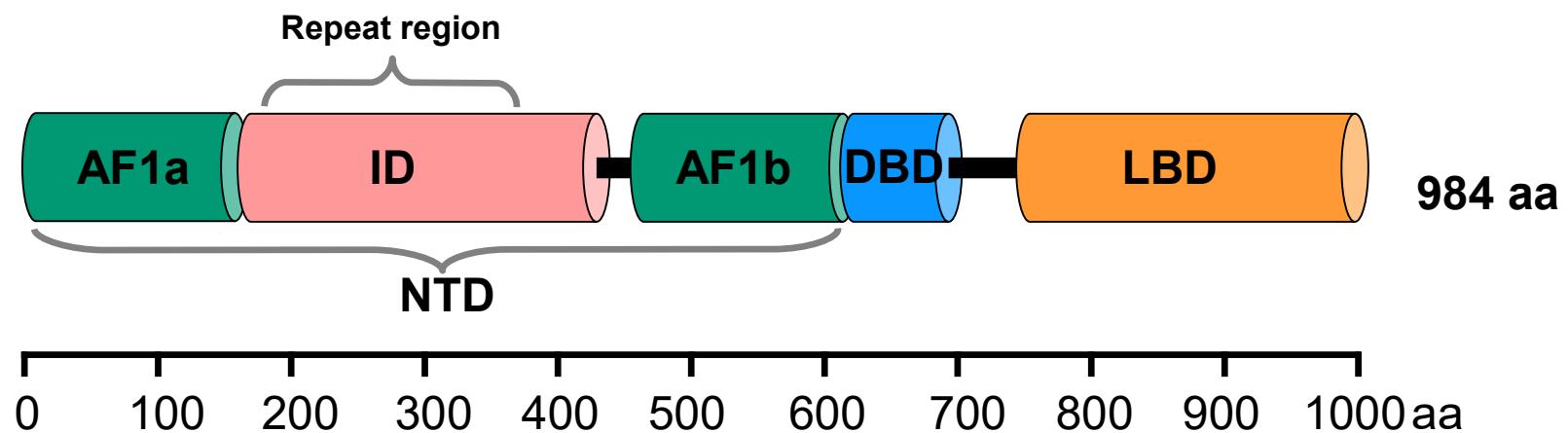
#F8

#F9

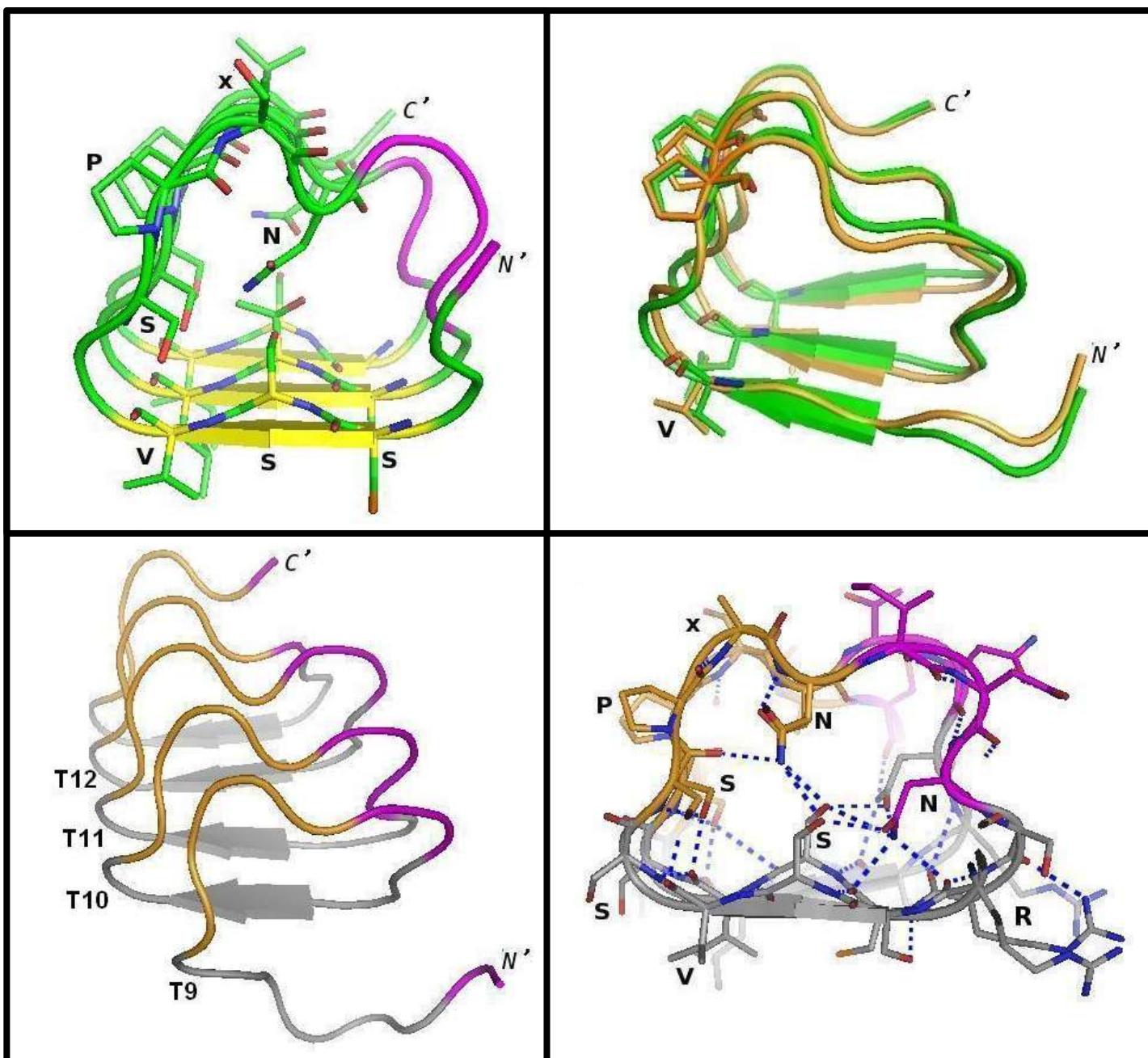
#F10

#F11

Mineralocorticoid receptor



Vlassi *et al.* (2013) *BMC Struct. Biol.*



Composition bias

Definition

14% proteins contains repeats (Marcotte et al, 1999)

1: Single amino acid repeats.

2: Longer imperfect tandem repeats.
Assemble in structure.

Definition CBRs

Perfect repeat: QQQQQQQQQQQQQQ

Imperfect: QQQQPQQQQQQQ

Amino acid type: DDDDDDEEEDEED

Compositionally biased regions (CBRs)

High frequency of one or two amino acids in a region.

Particular case of low complexity region

Function CBRs

Conservation => Function

Length, amino acid type not necessarily conserved

Frequency: 1 in 3 proteins contains a compositionally biased region (Wootton, 1994), ~11% conserved (Sim and Creamer, 2004)

Function CBRs

Conservation => Function

Length, amino acid type not necessarily conserved

Functions:

Passive: linkers

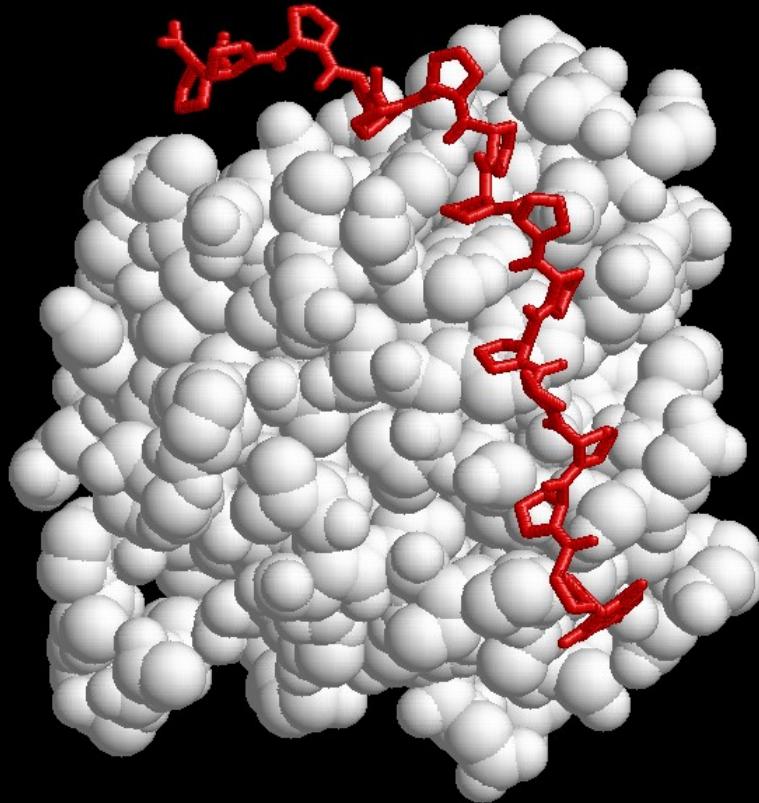
Active: binding, mediate protein interaction, structural integrity

(Sim and Creamer, 2004)

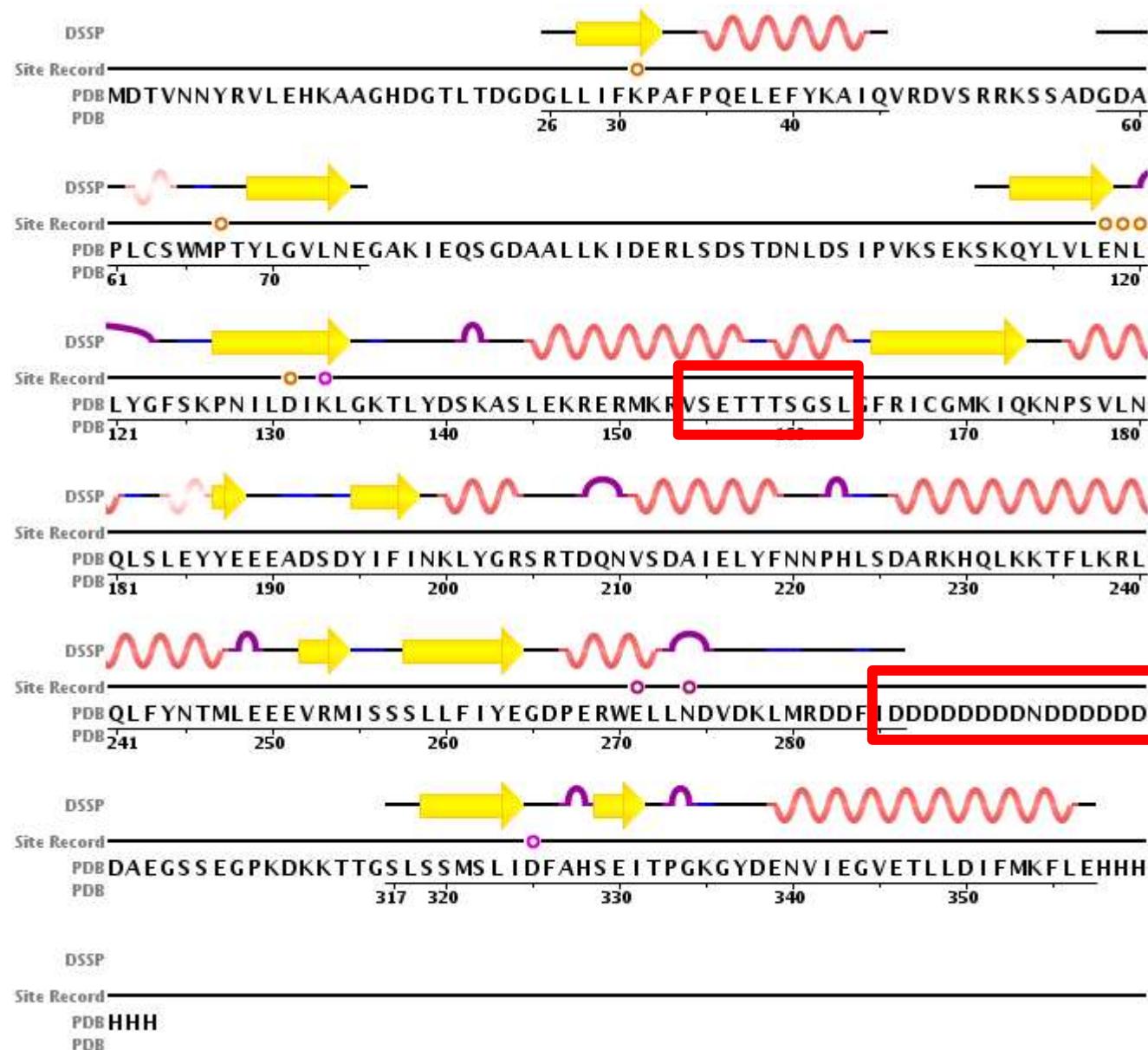
Structure of CBRs

Often variable or flexible: do not easily
crystallize

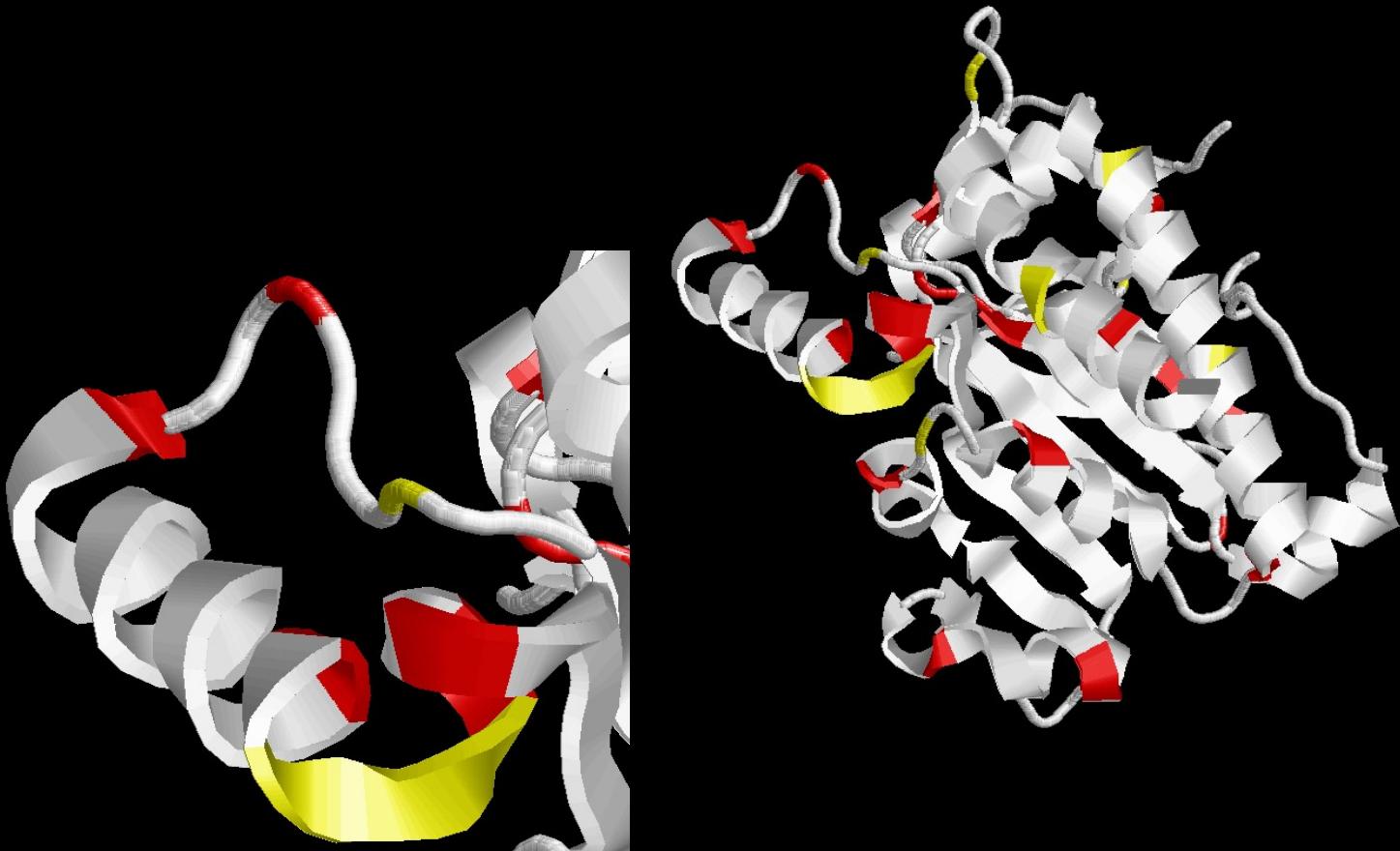
1CJF: profilin bound to polyP



2IF8: Inositol Phosphate Multikinase Ipk2

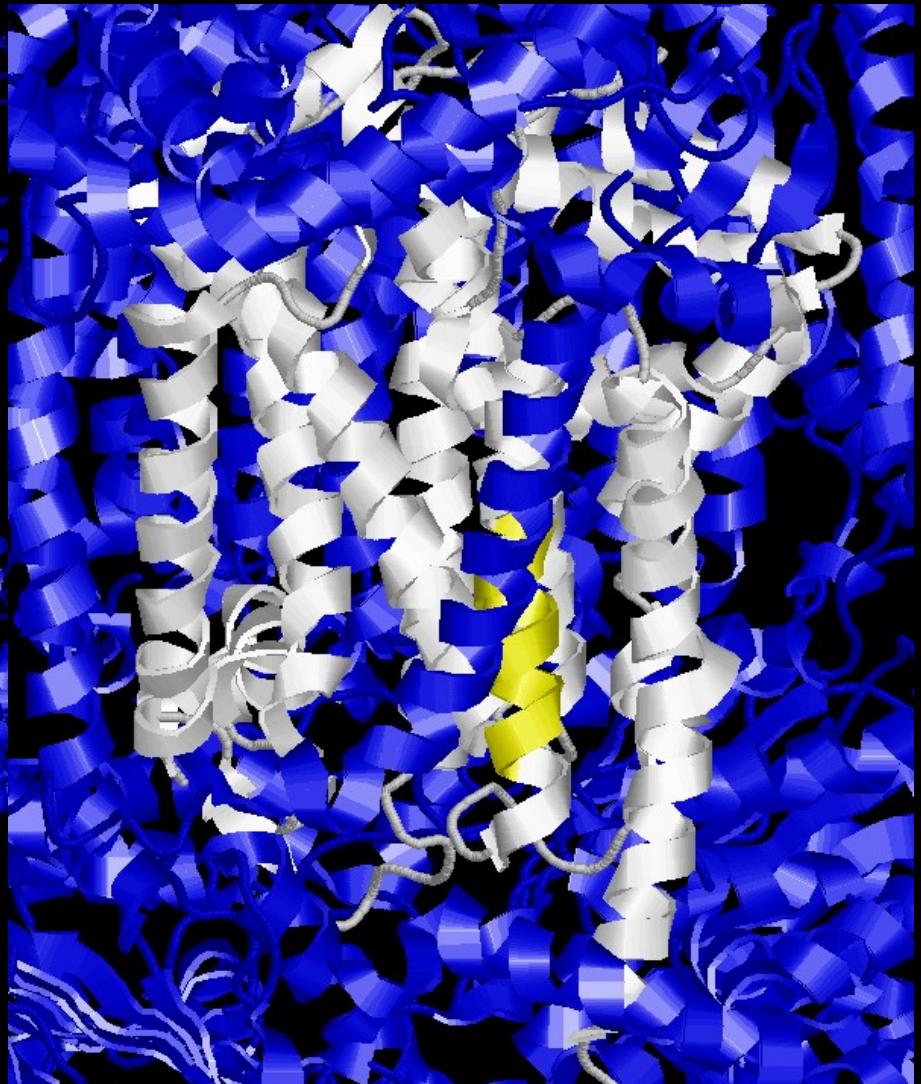
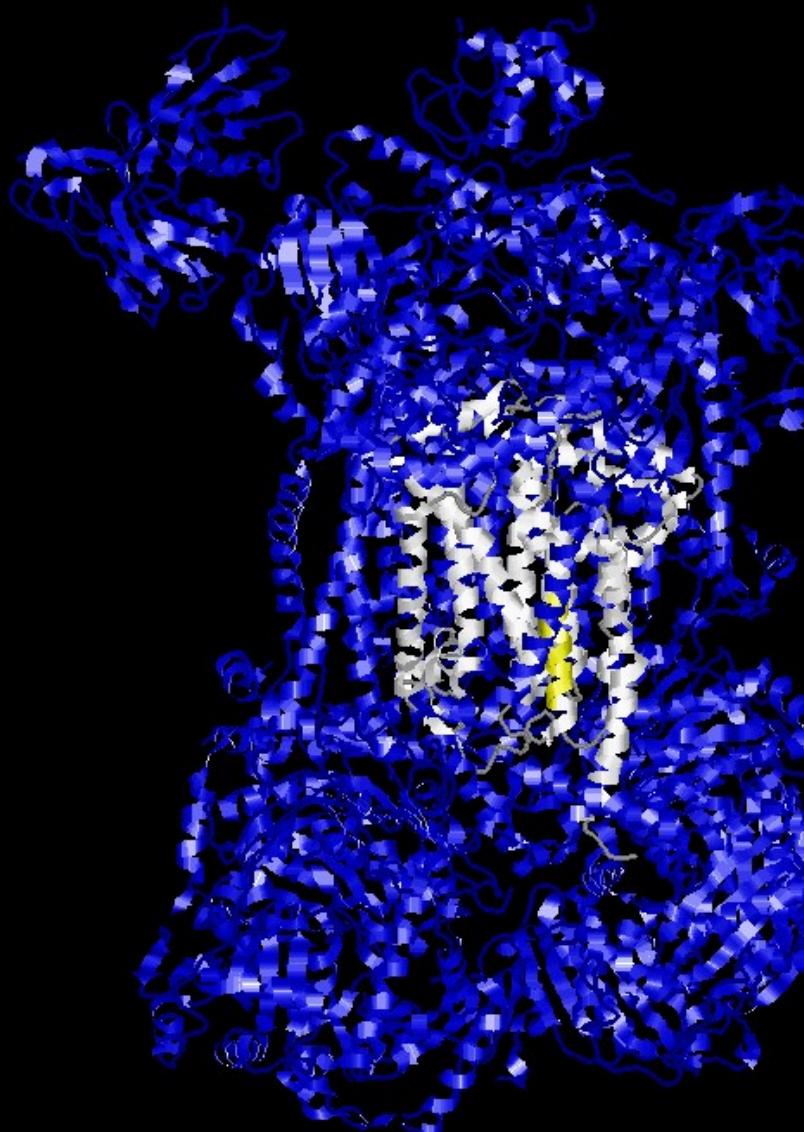


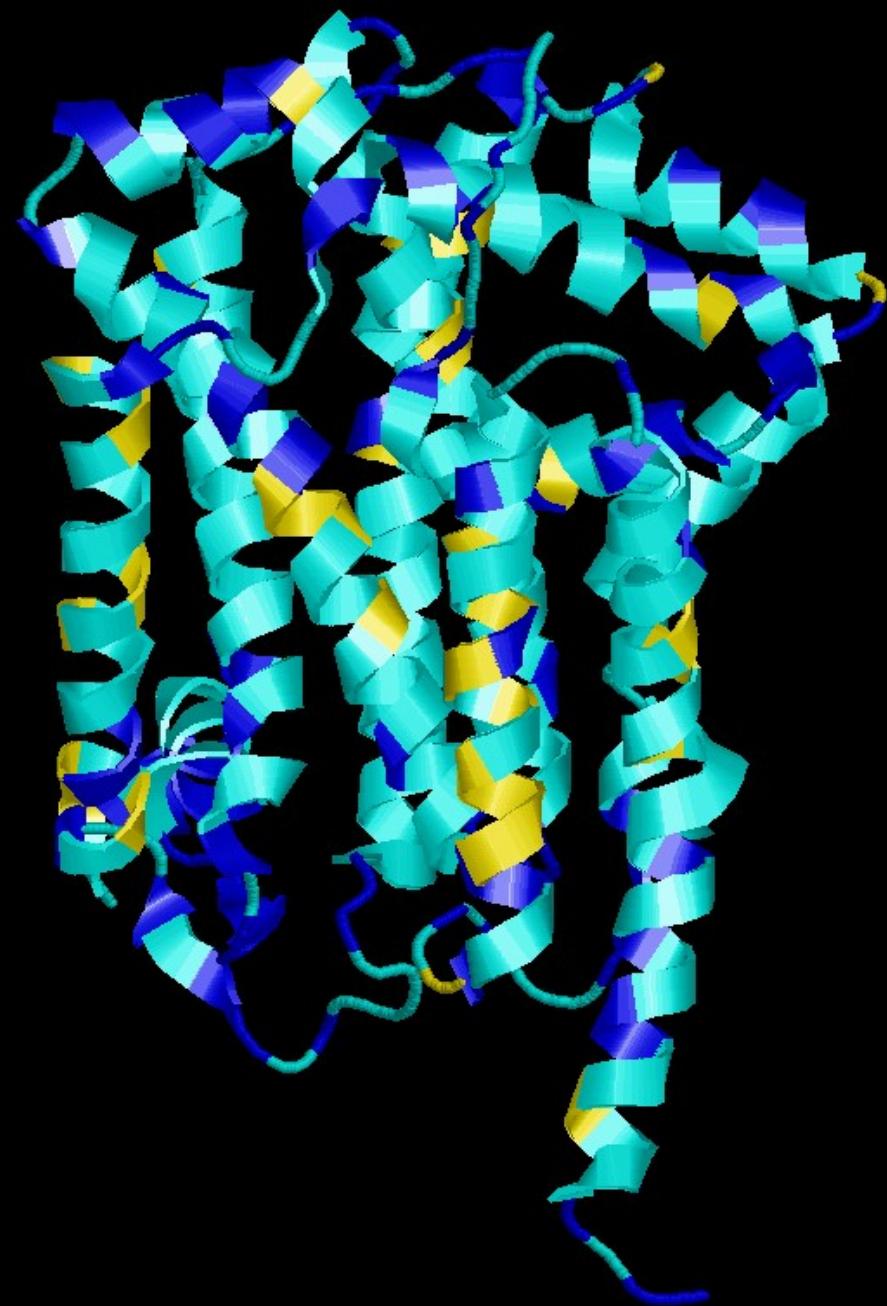
2IF8: Inositol Phosphate Multikinase Ipk2



RVSETTTSGSL

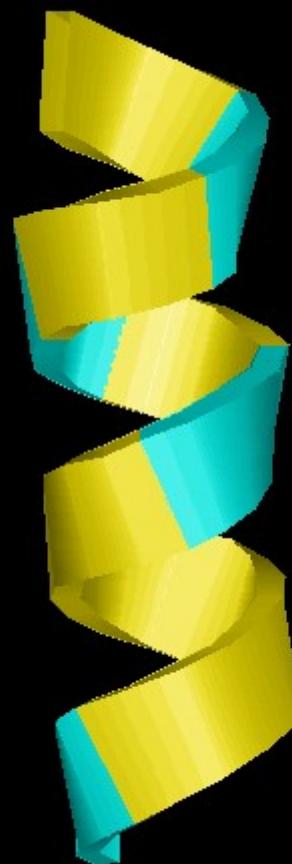
2CX5: mitochondrial
cytochrome c
B subunit N-terminal





2CX5: mitochondrial
cytochrome c
B subunit N-terminal

EFFECTIVENESS



Amino acid repeats

Distribution is not random:

Eukaryota:

Most common: poly-Q, poly-N, poly-A, poly-S, poly-G

Prokaryota:

Most common: poly-S, poly-G, poly-A, poly-P

Relatively rare: poly-Q, poly-N

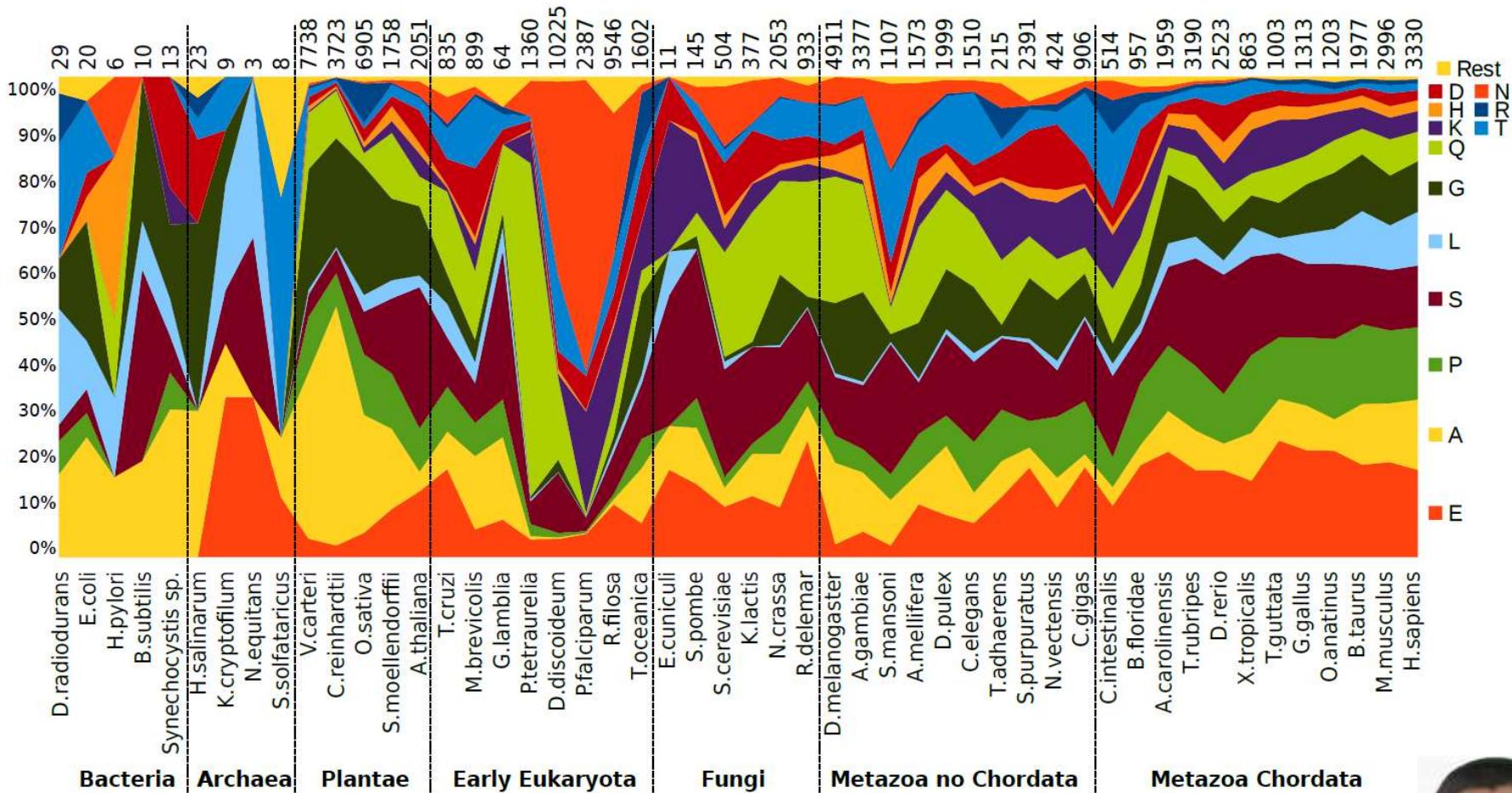
Very rare or absent in both eukaryota and prokaryota:

Poly-I, Poly-M, Poly-W, Poly-C, Poly-Y

Toxicity of long stretches of hydrophobic residues.

(Faux et al 2005)

Amino acid repeats



Mier et al. (2017) Proteins

Pablo
Mier



Filtering out CBRs

Normally filtered out as low complexity region: they give spurious BLAST hits

QQQQQQQQQQ

||||||| | | |

QQQQQQQQQQ 10/10 id

IDENTITIES

||||||| | | |

IDENTITIES 10/10 id

Filtering out CBRs

Normally filtered out as low complexity region: they give spurious BLAST hits

QQQQQQQQQQ

|||||||

QQQQQQQQQQ Shuffle: 10/10 id

IDENTITIES

|||||||

IDENTITIES 10/10 id

Filtering out CBRs

Normally filtered out as low complexity region: they give spurious BLAST hits

QQQQQQQQQQ

|||||||

QQQQQQQQQQ Shuffle: 10/10 id

IDENTITIES

| |

SIINDIETTE Shuffle: 2/10 id

Filtering out CBRs

Option for pre-BLAST treatment

SEG algorithm:

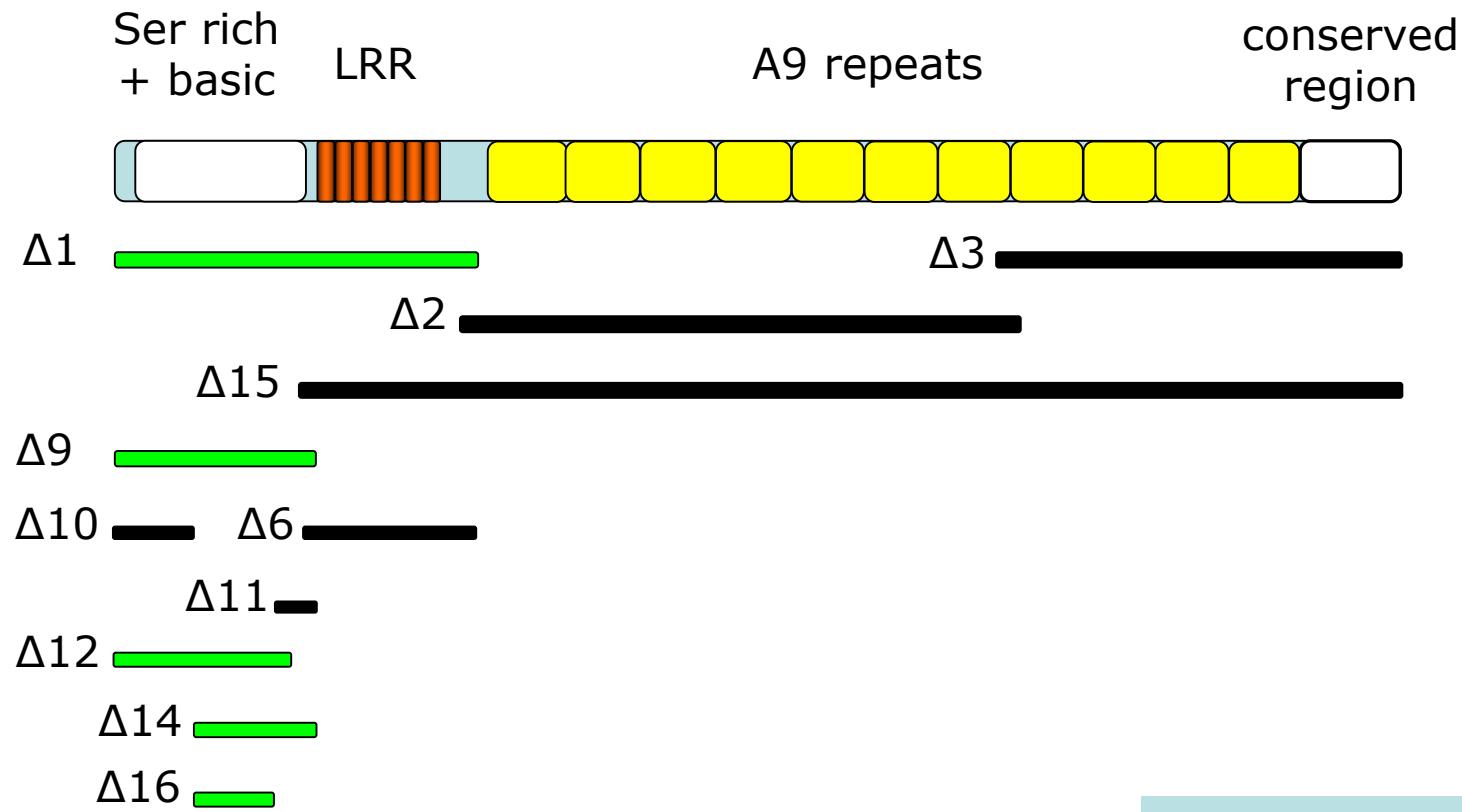
- 1) Identify sequence regions with low information content over a sequence window
- 2) Merge neighbouring regions

Eliminates hits against common acidic-, basic- or proline-rich regions

(Wootton and Federhen, 1993)

AIR9

(1708 aa)



Microtubule localization of Δx -GFP

Buschmann, et al (2006).
Current Biology.
Buschmann, et al (2007).
Plant Signaling & Behavior

Homorepeats are frequent but difficult to characterize

Pablo
Mier



e.g. polyQ:

MATLEKLMKAFESLKSFQQQQQQQQQQQQQQQQQQQQQQQQQQQQQQPPPPPPPPPPPPPQLPQP

- 10% of human proteins have homorepeats
- lack sequence conservation
- not possible to predict function by homology

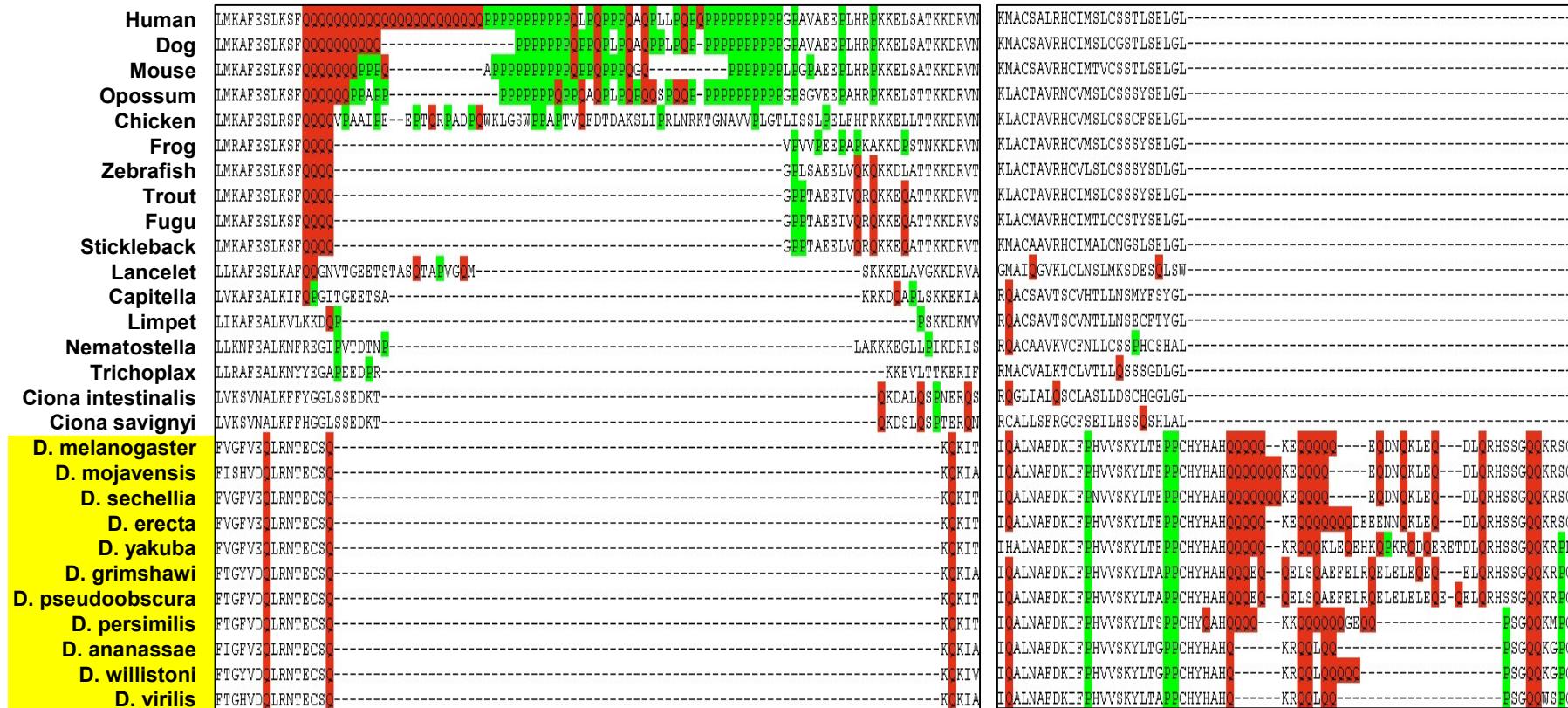
Homorepeats need to be studied in context

Martin
Schaefer

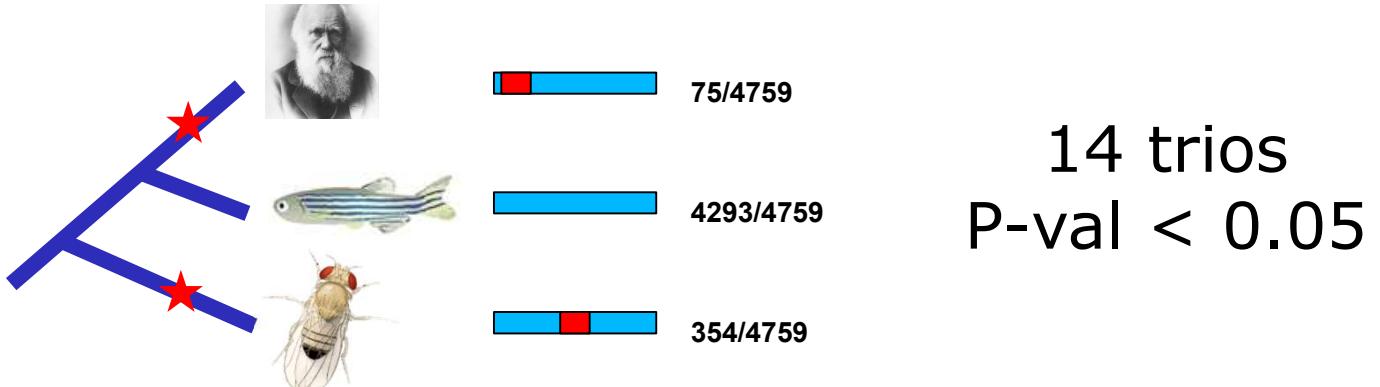
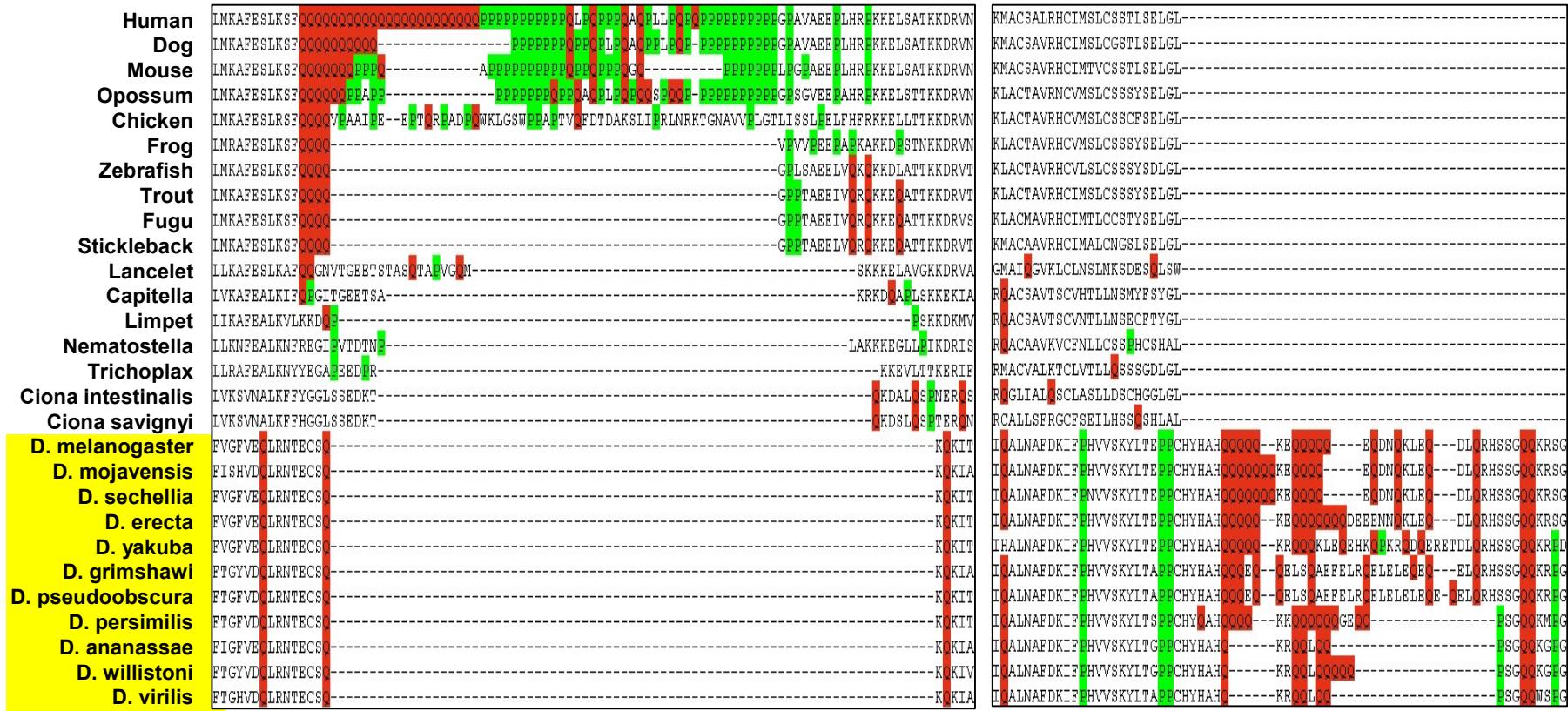


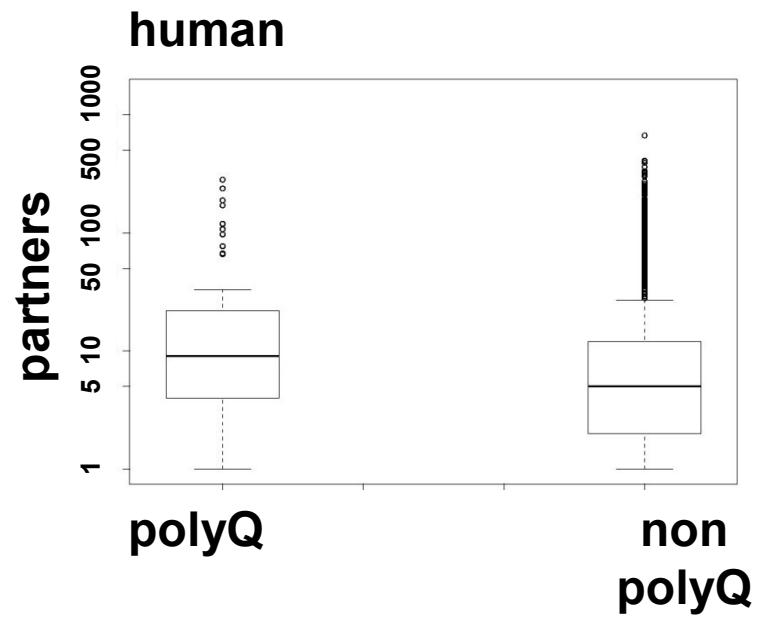
Function of polyQ

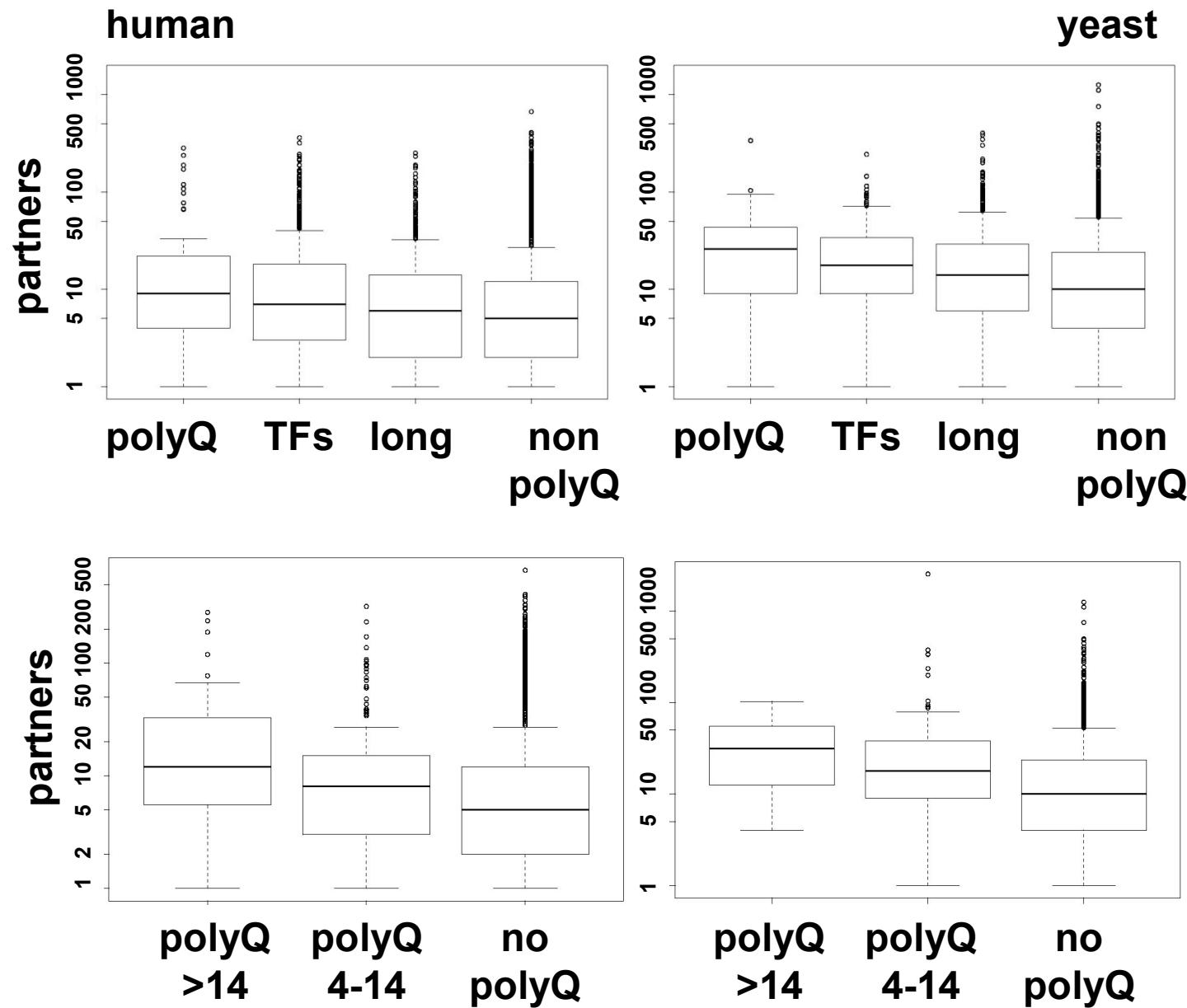
polyQ in Huntingtin



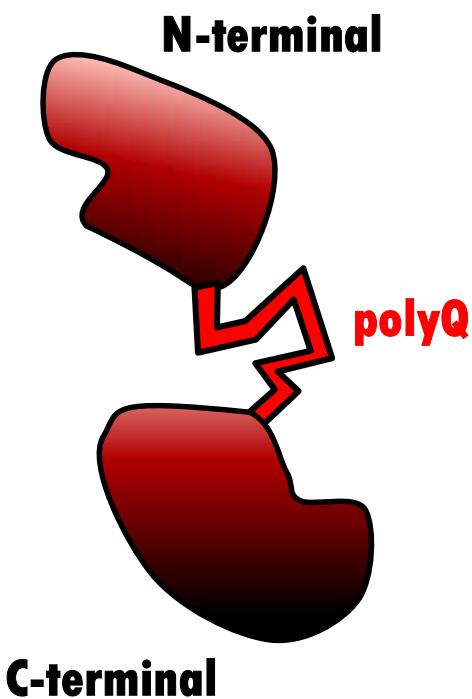
Schaefer et al (2012) *Nucleic Acids Res.*



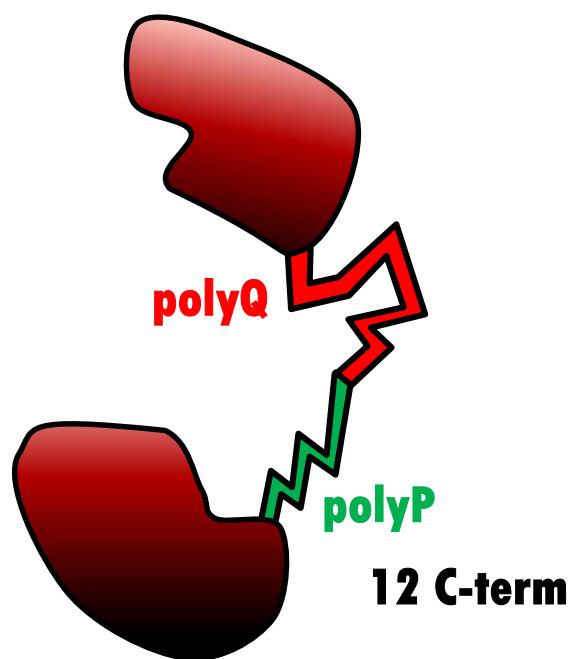




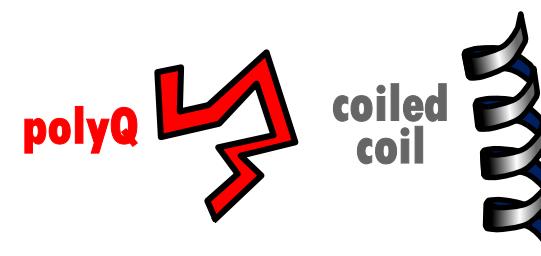
**86 human
polyQ
proteins**



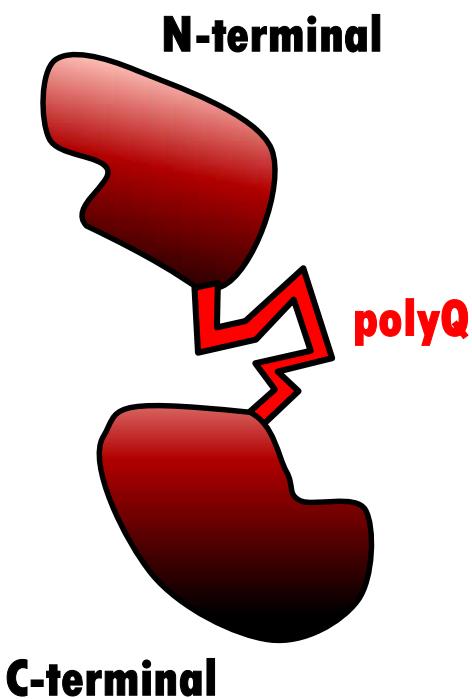
**13 polyQ
proteins with
near polyP**



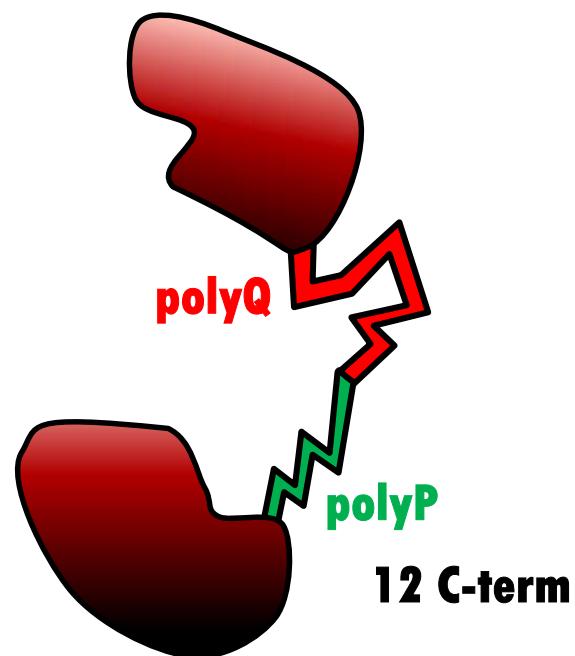
**109 polyQ regions
54 overlap/near
coiled coil**



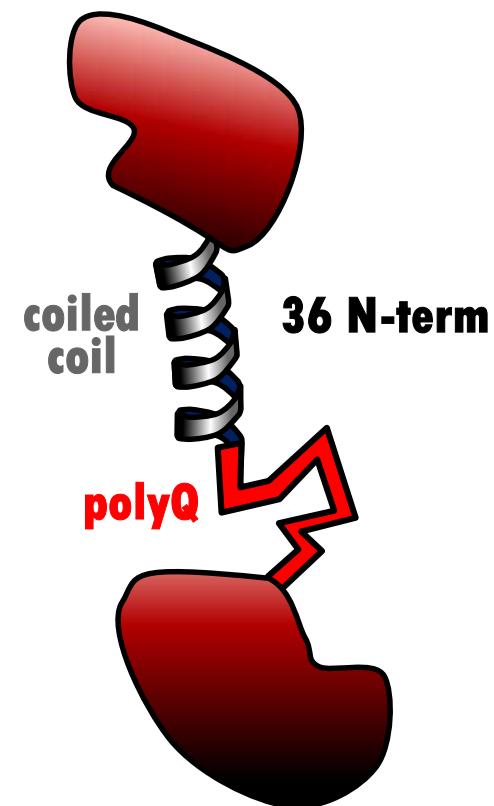
**86 human
polyQ
proteins**



**13 polyQ
proteins with
near polyP**

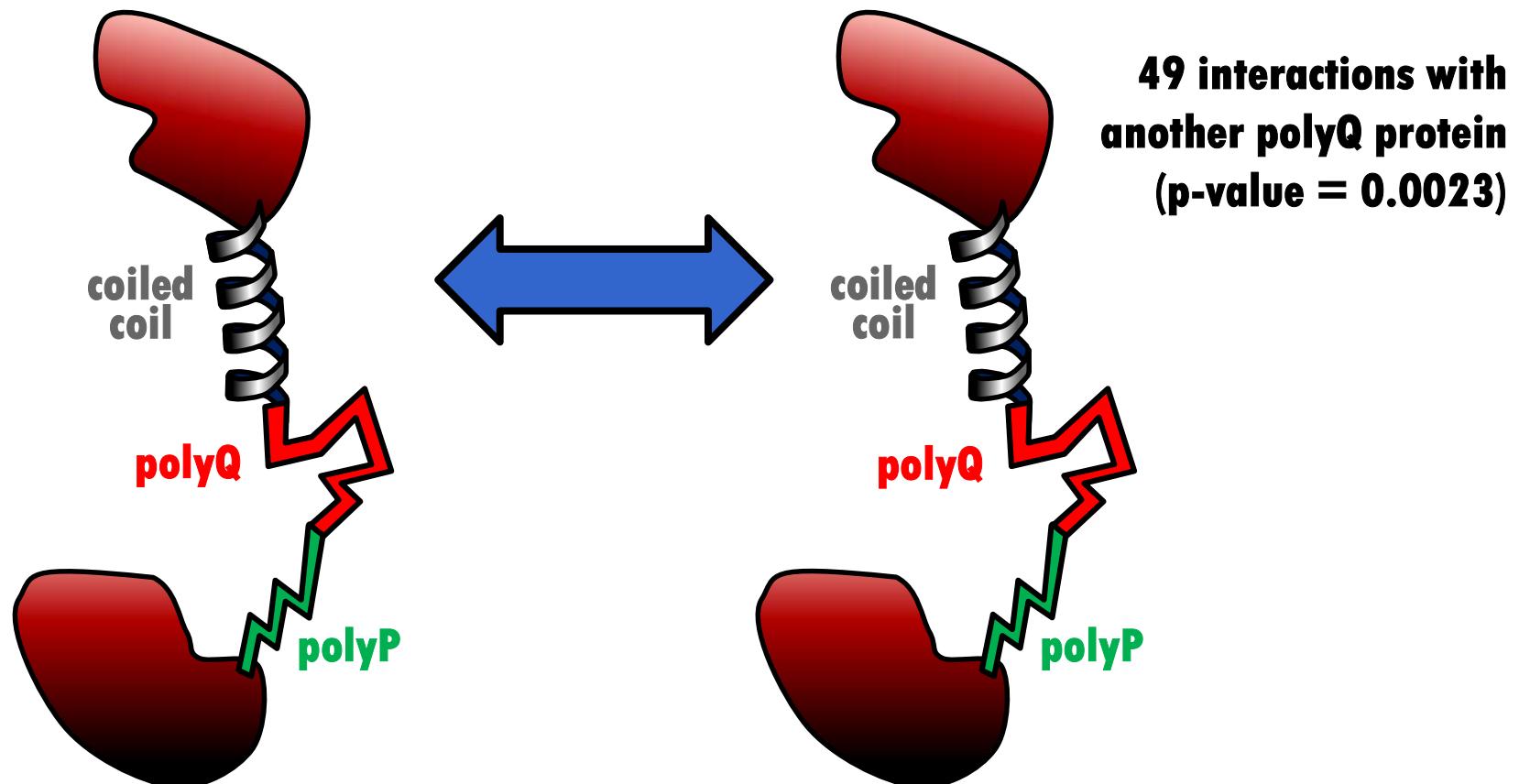


**40 human
polyQ/coiled-coil
proteins
(no polyP)**



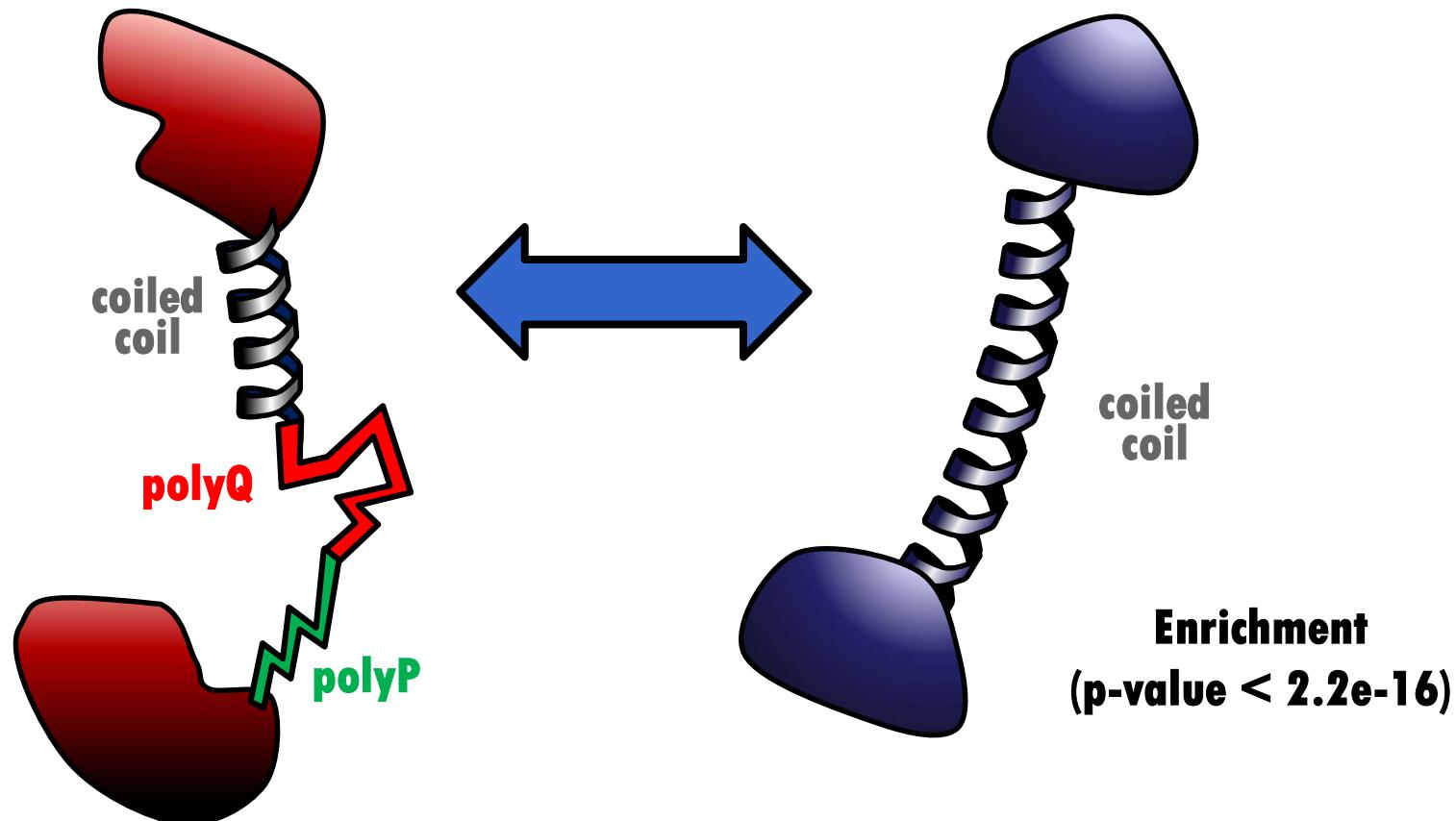
**86 human
polyQ
proteins**

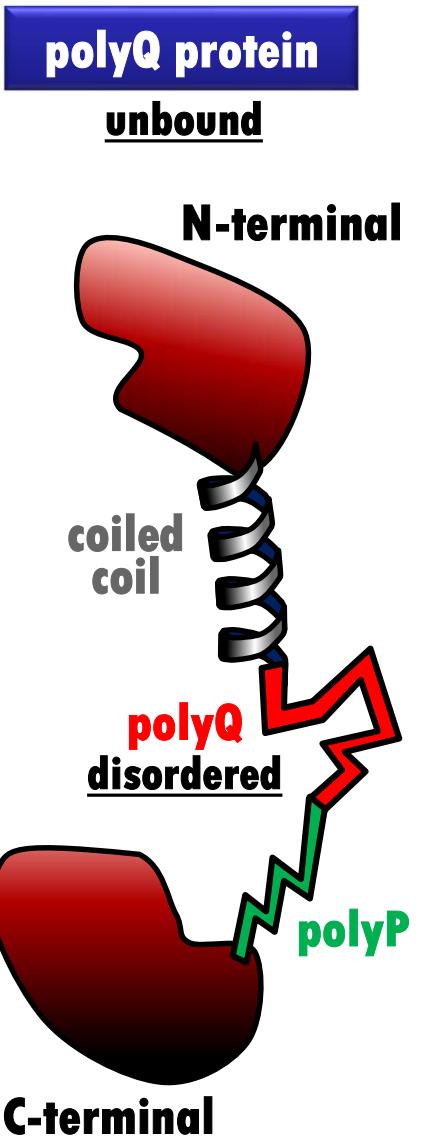
interacting proteins

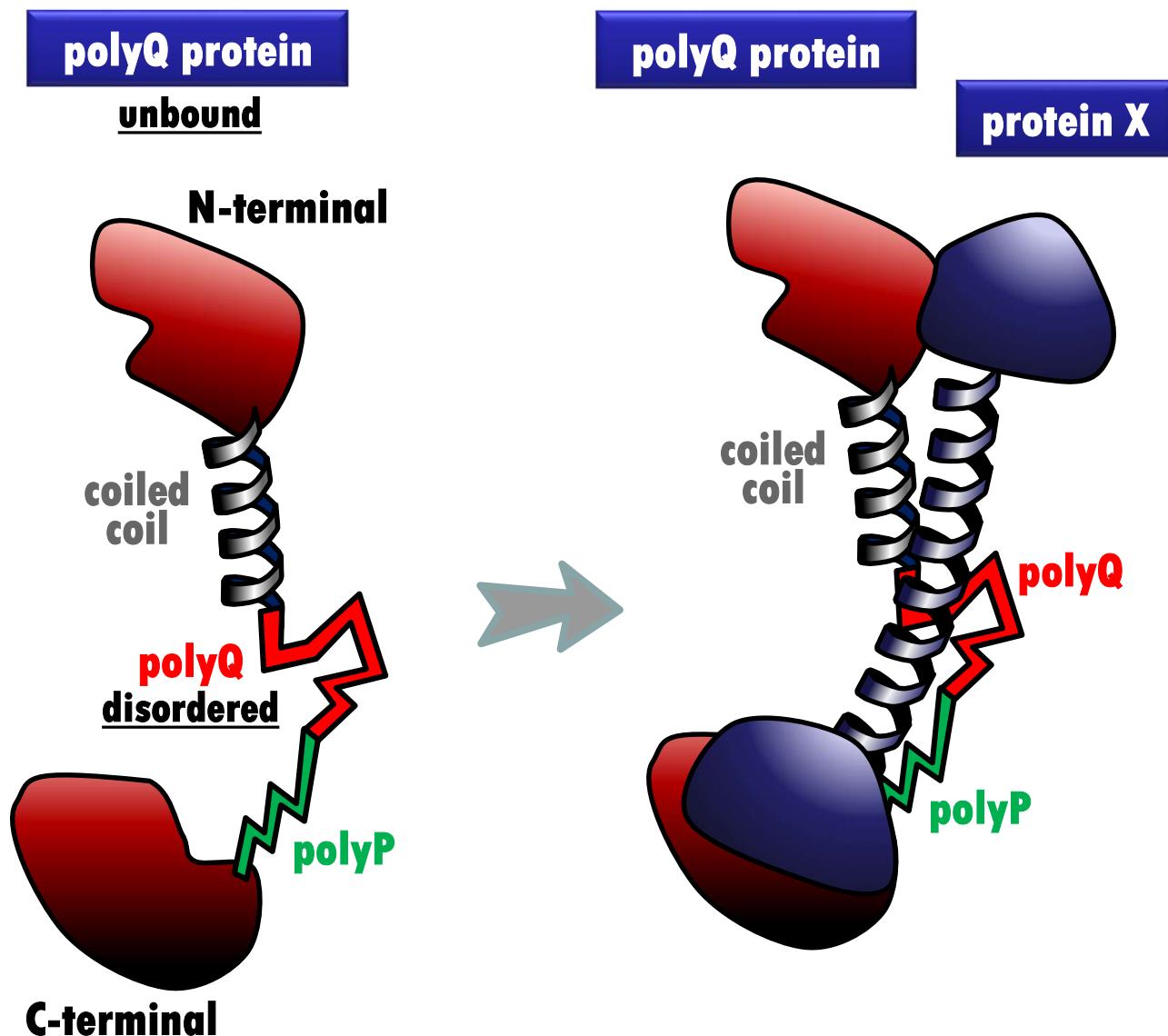


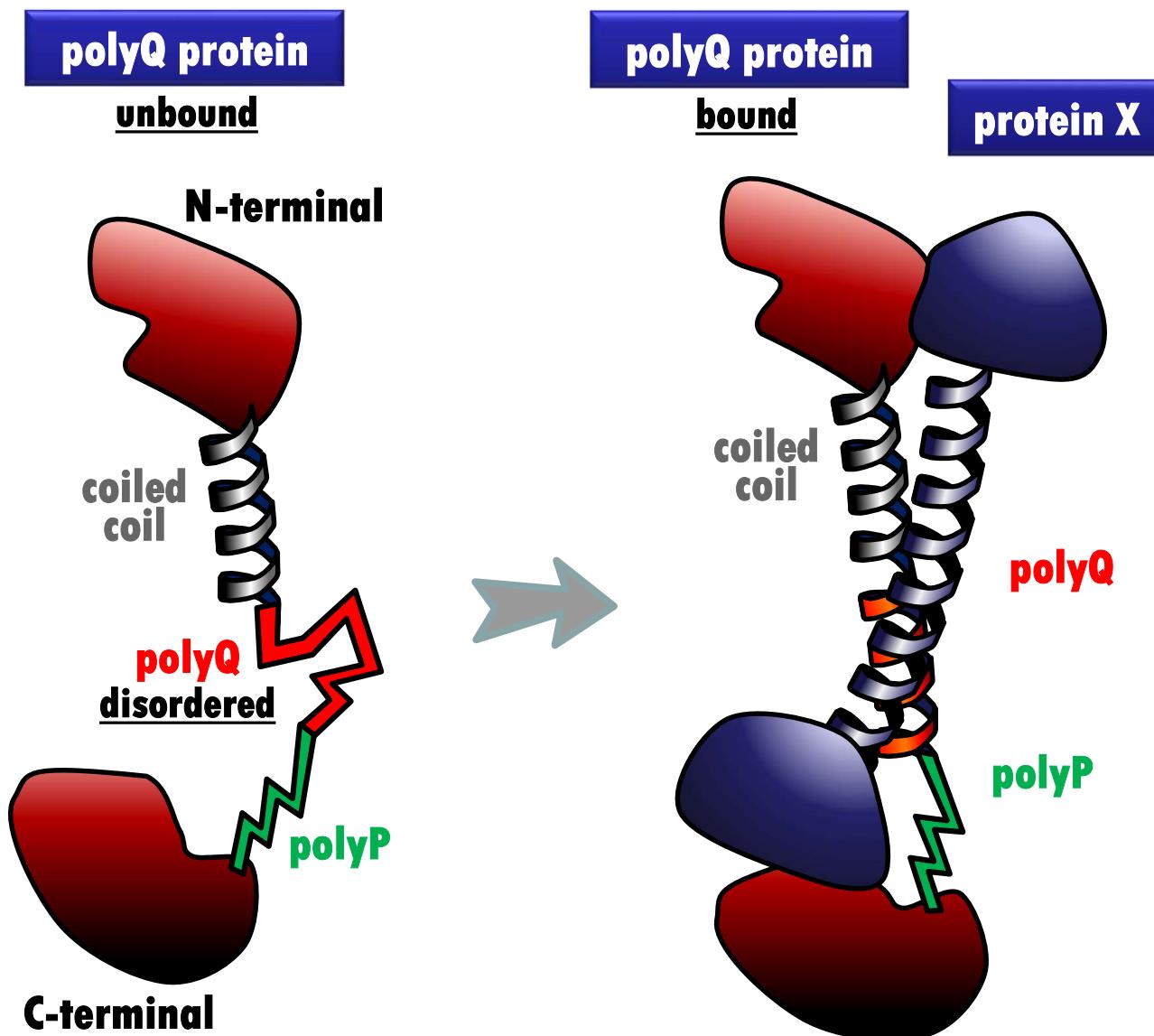
**86 human
polyQ
proteins**

**Non-polyQ
interacting proteins**



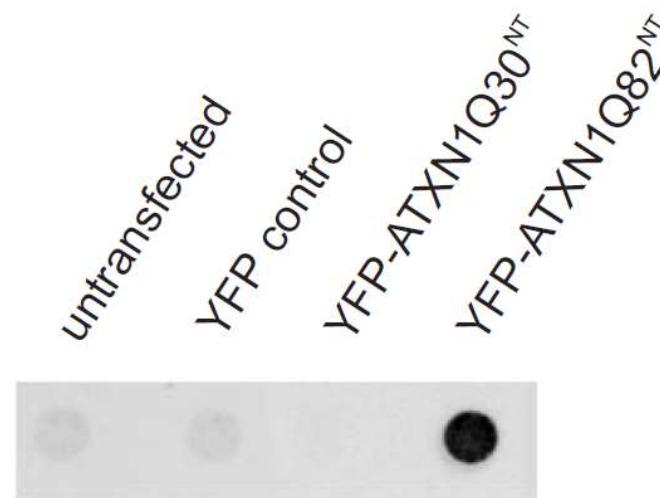
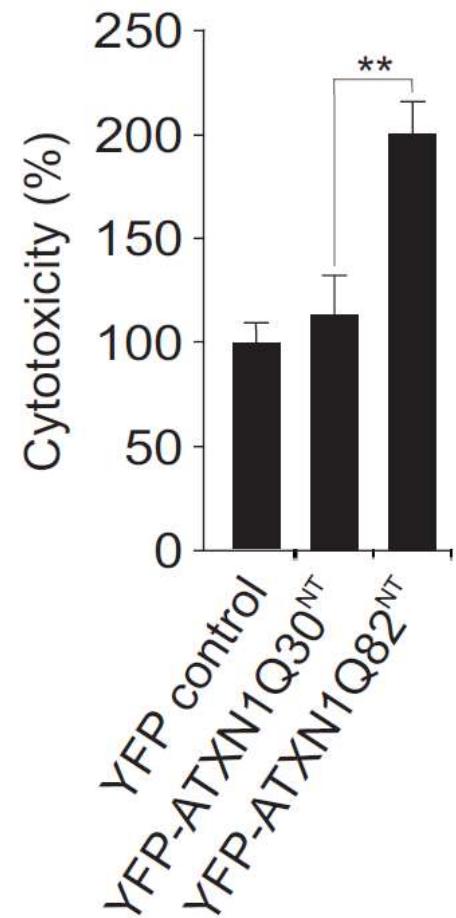






ATXN1Q82^{NT} is toxic

ATXN1Q82^{NT} aggregates



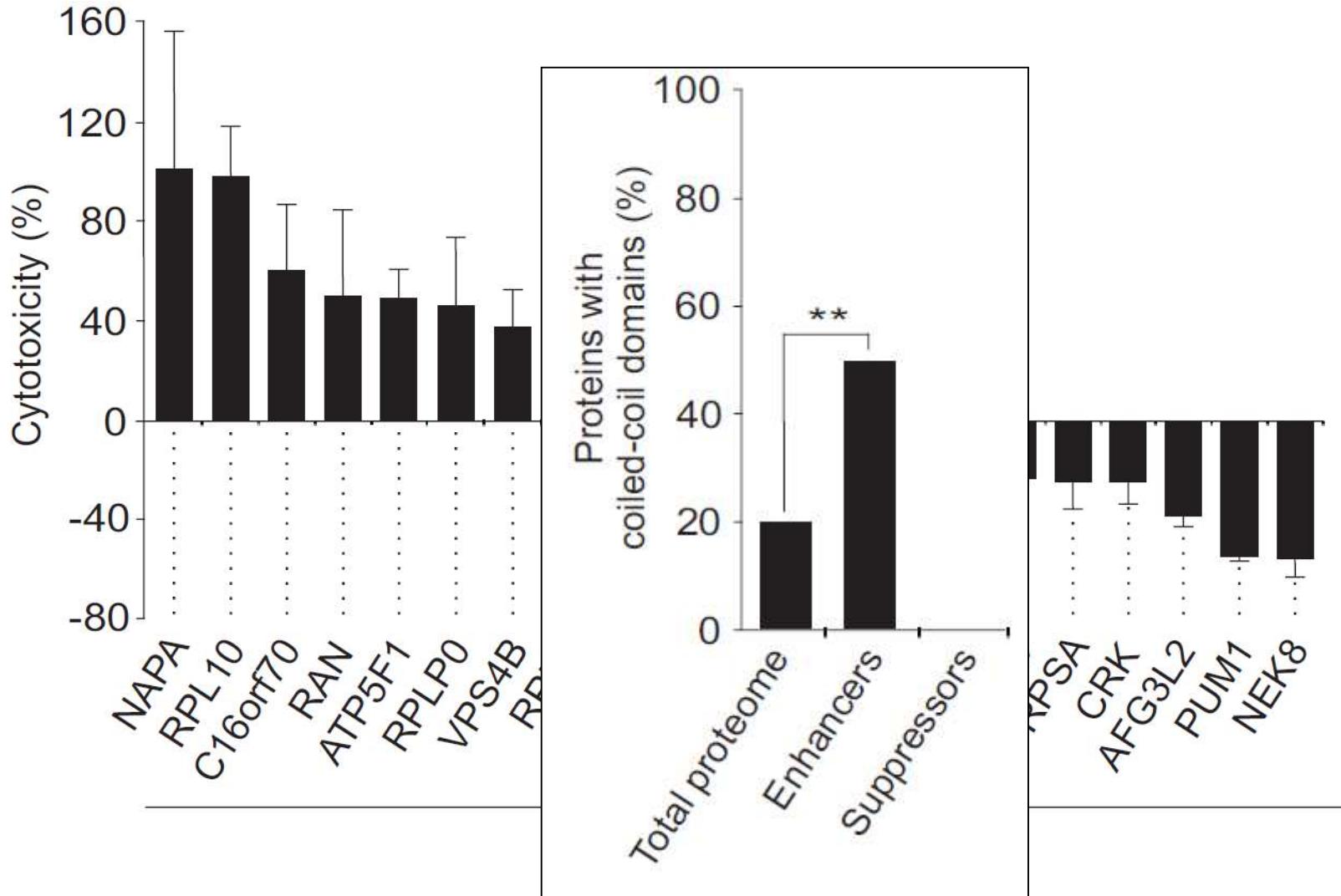
Spyros
Petrakis

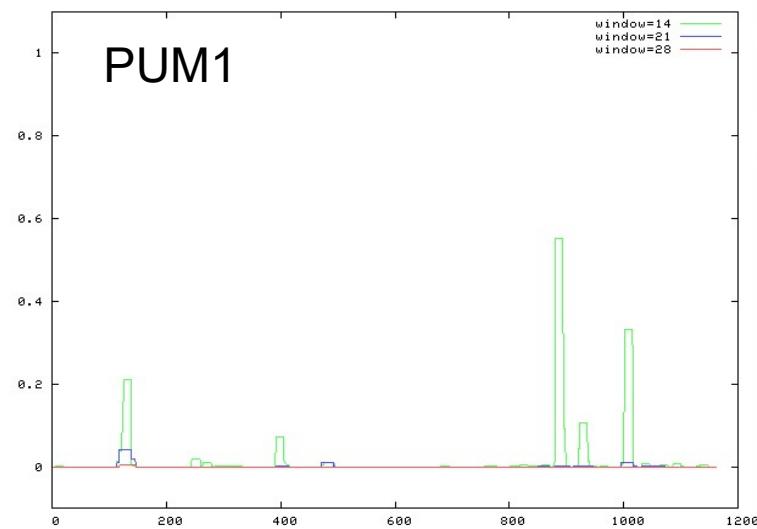
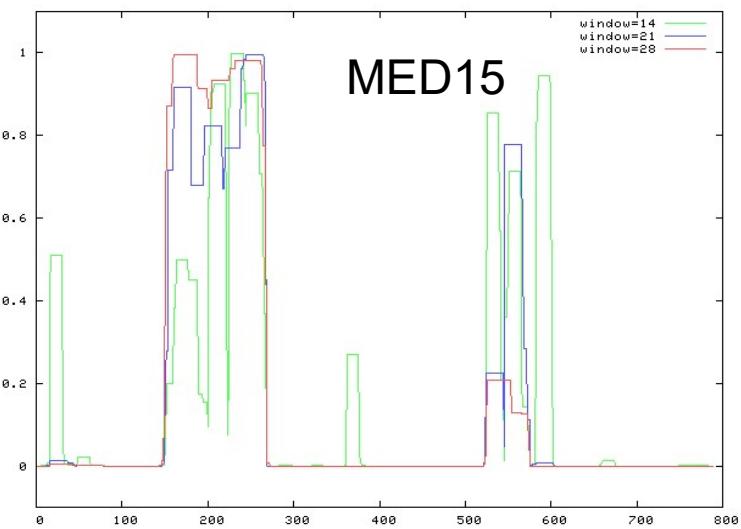
Erich
Wanker

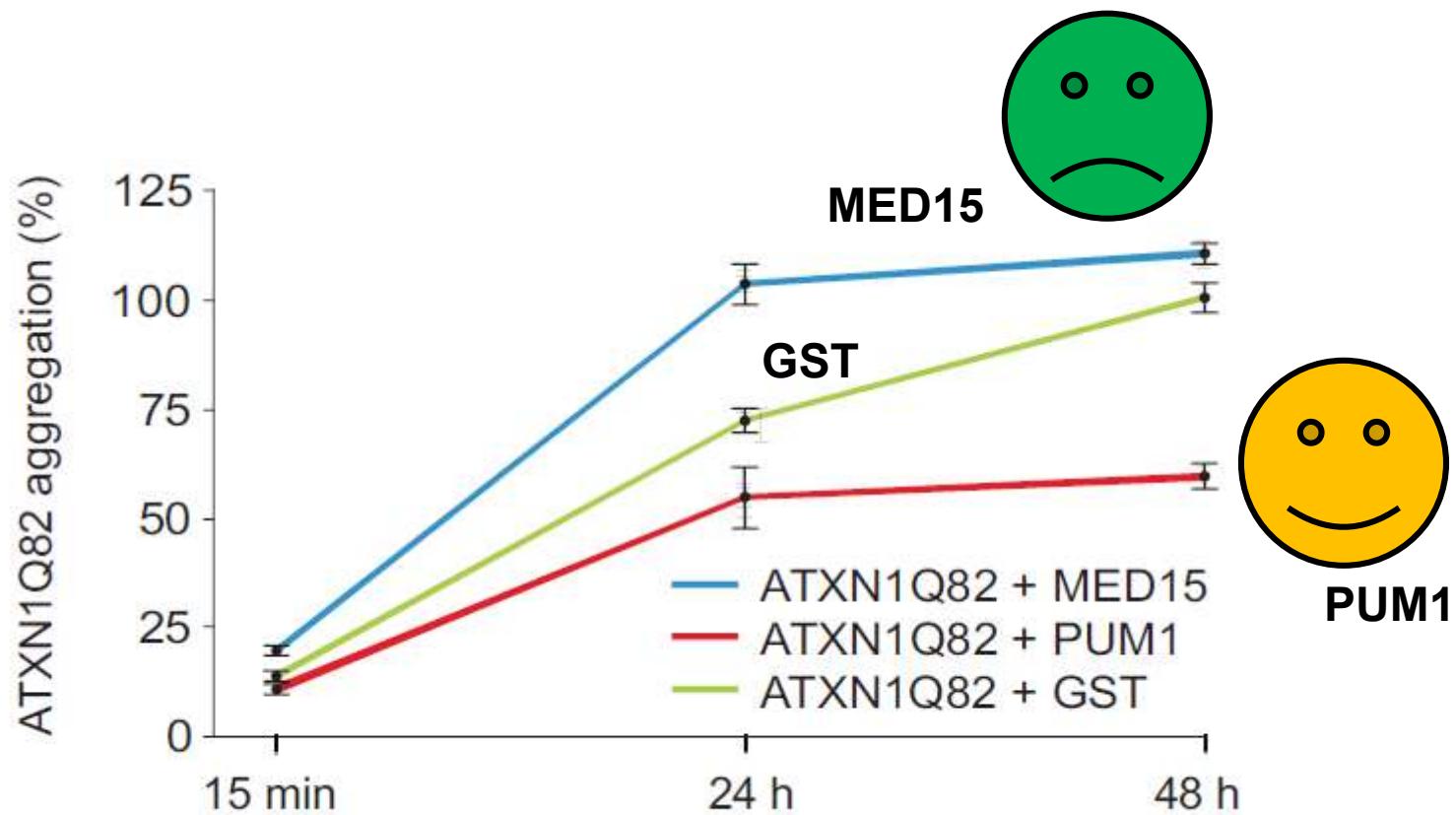


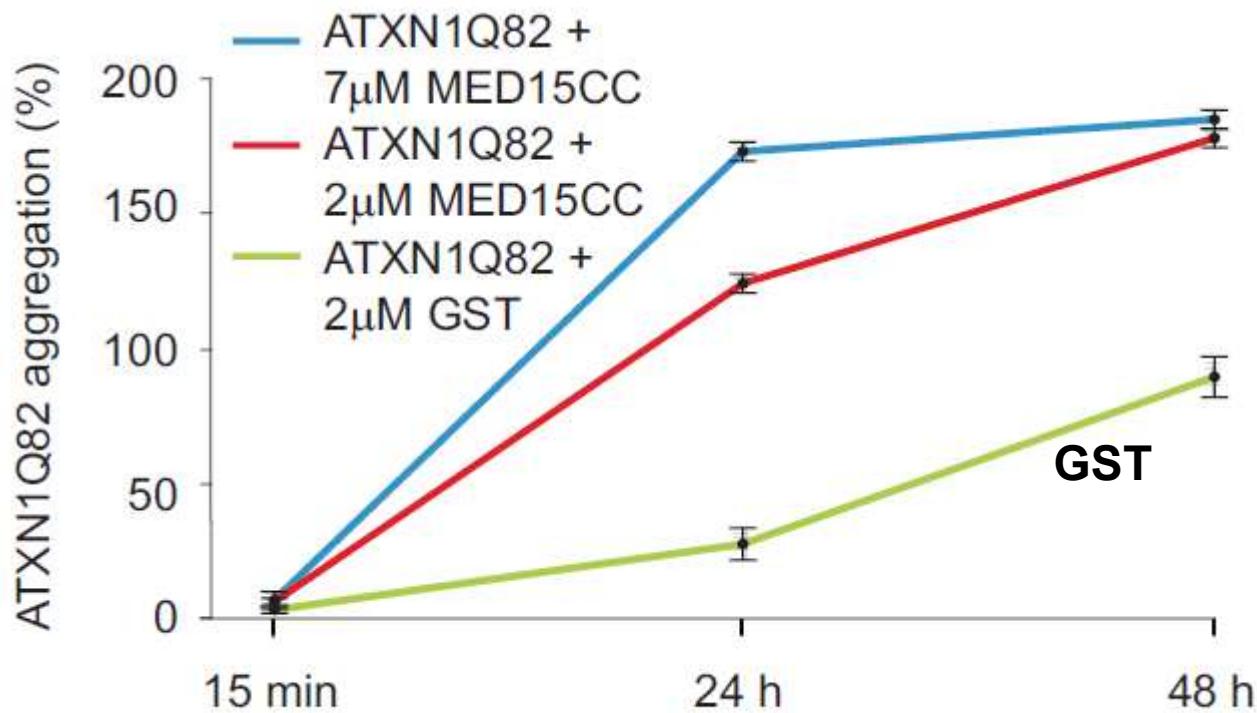
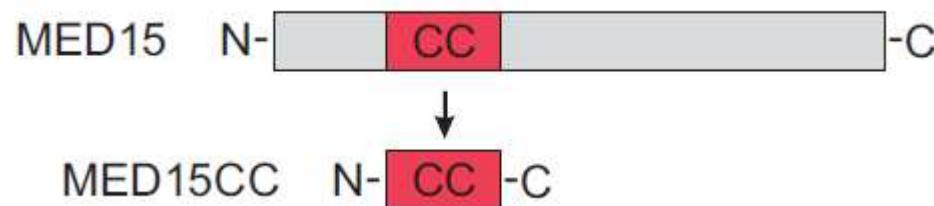
Petrakis et al. (2012) PLoS Genetics

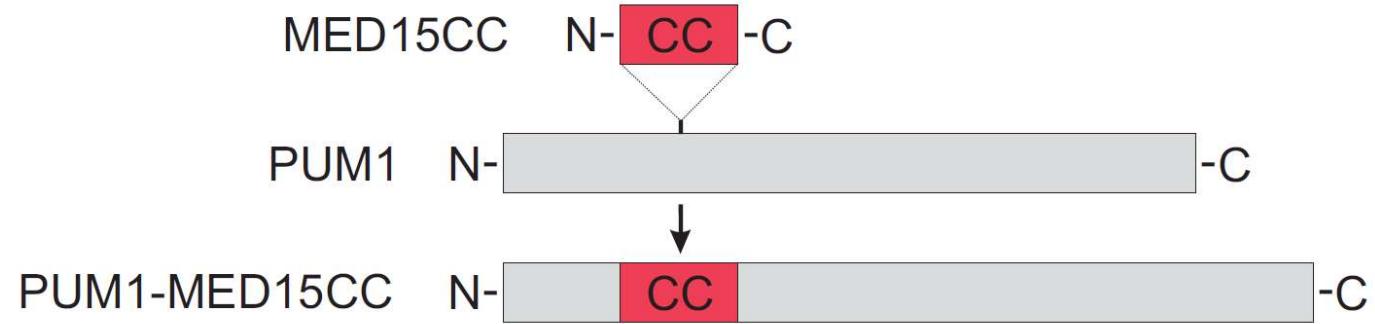
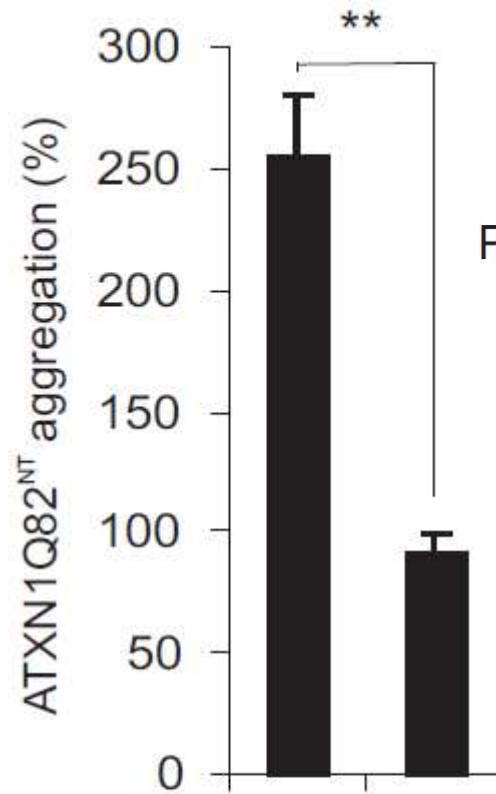
interactors that change ATXN1Q82^{NT} toxicity

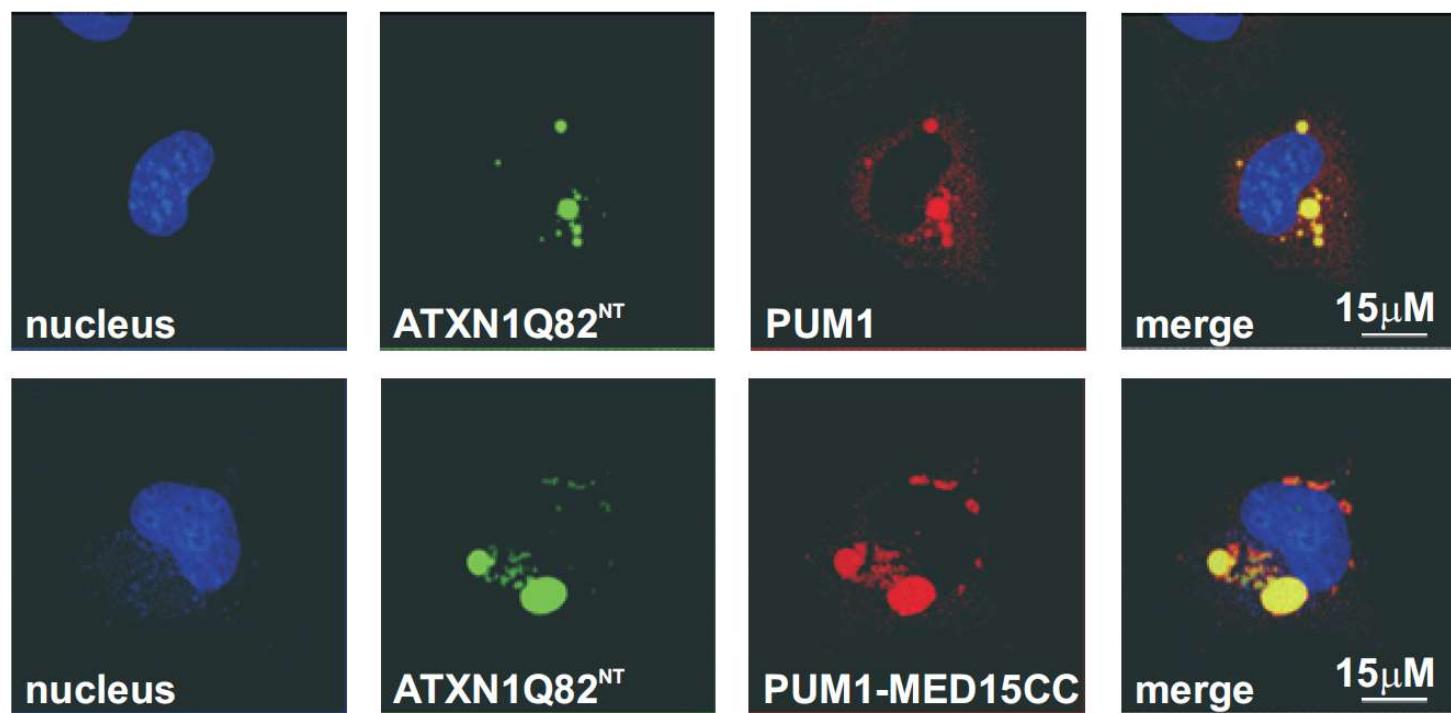
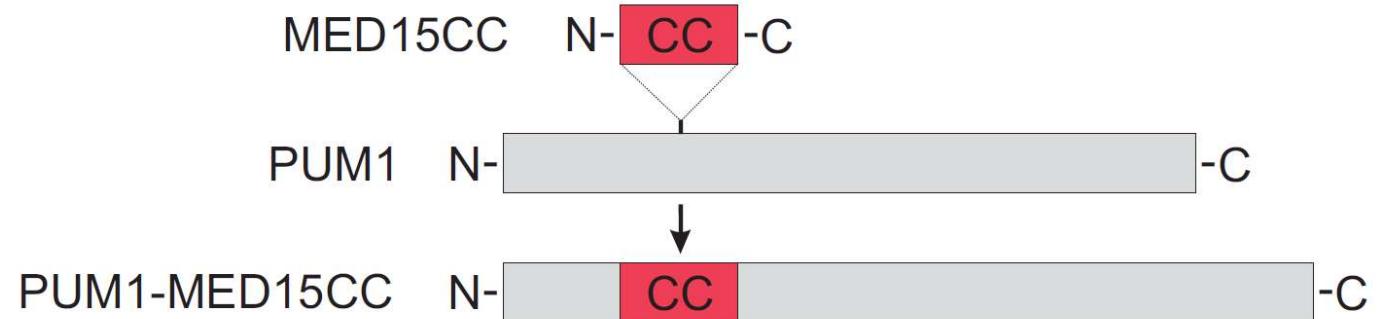


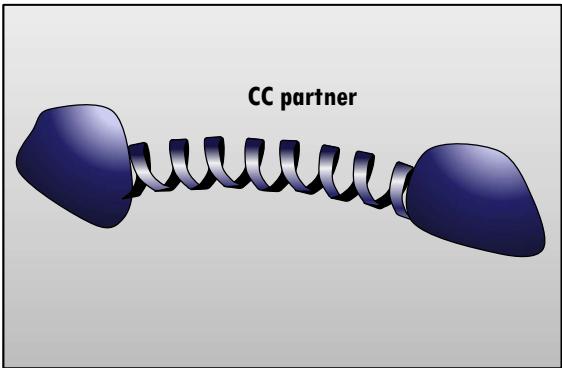
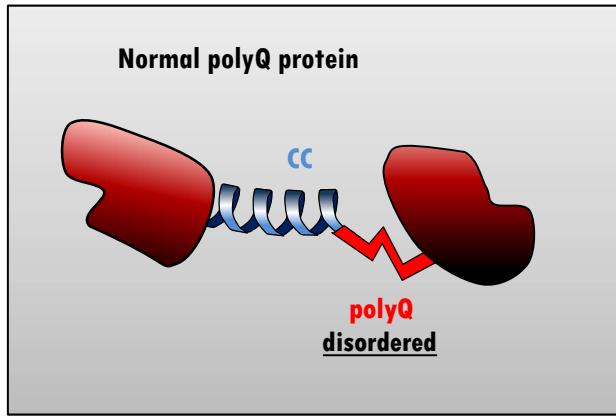


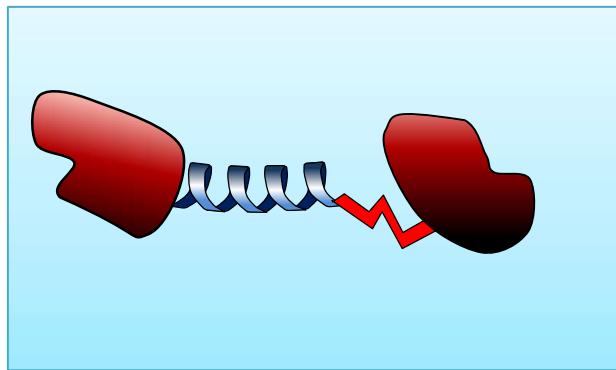
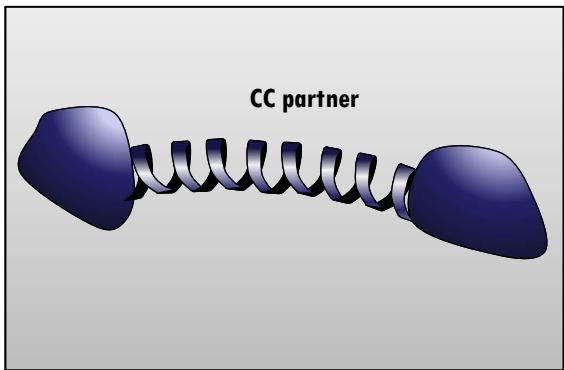
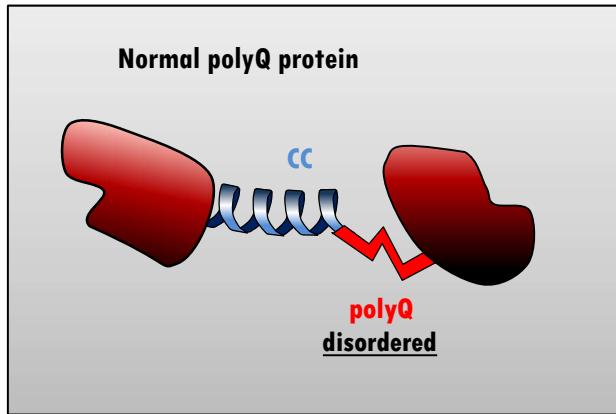


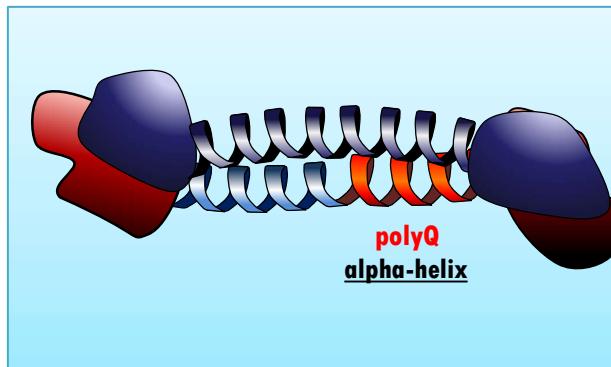
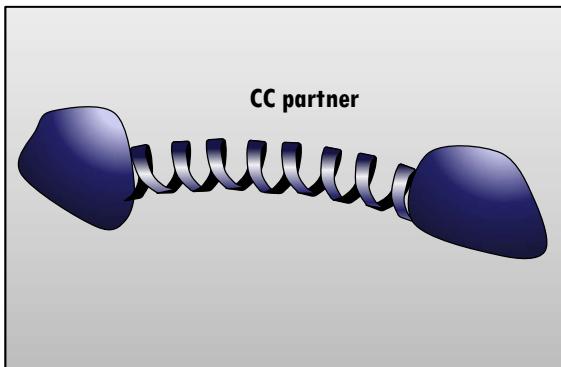
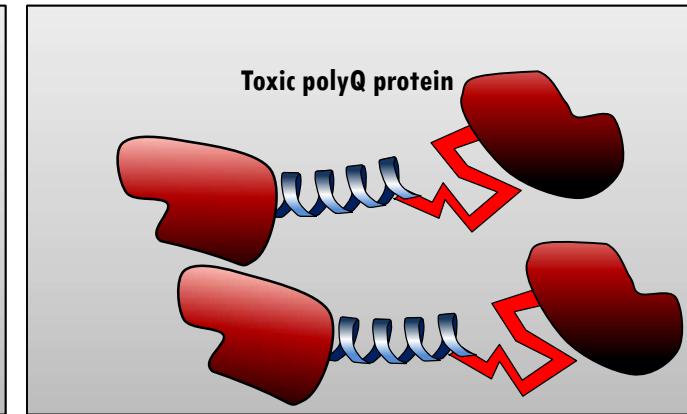
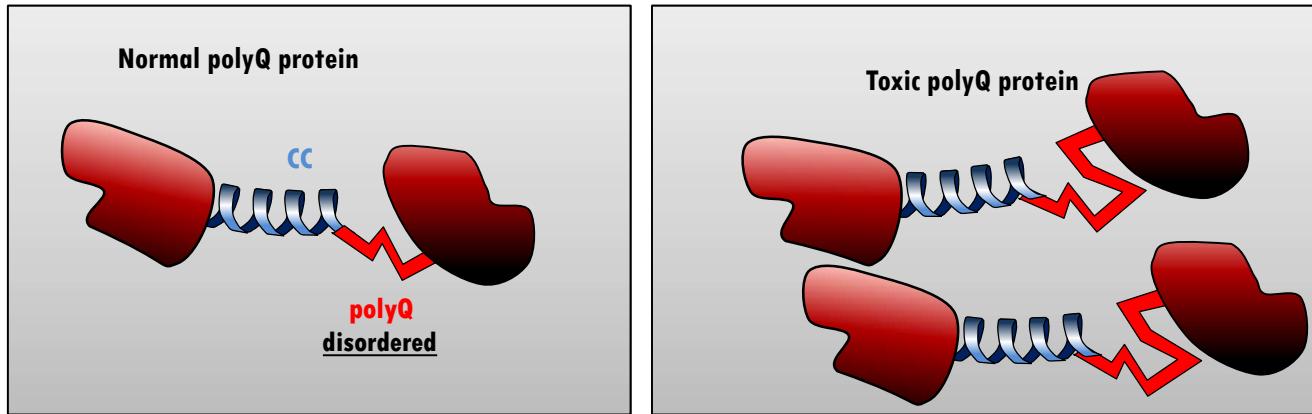


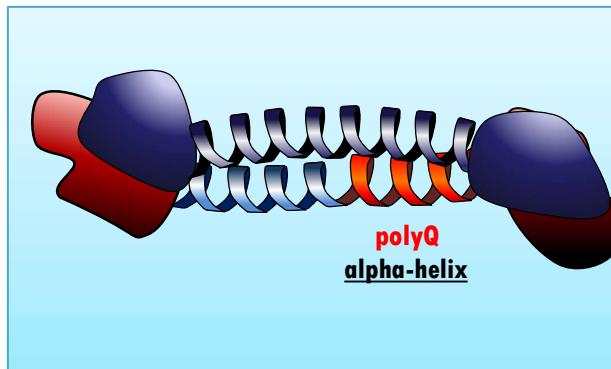
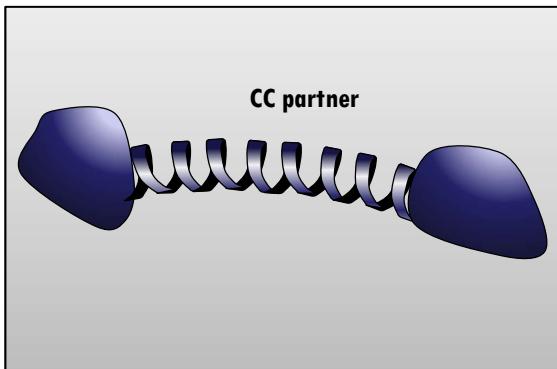
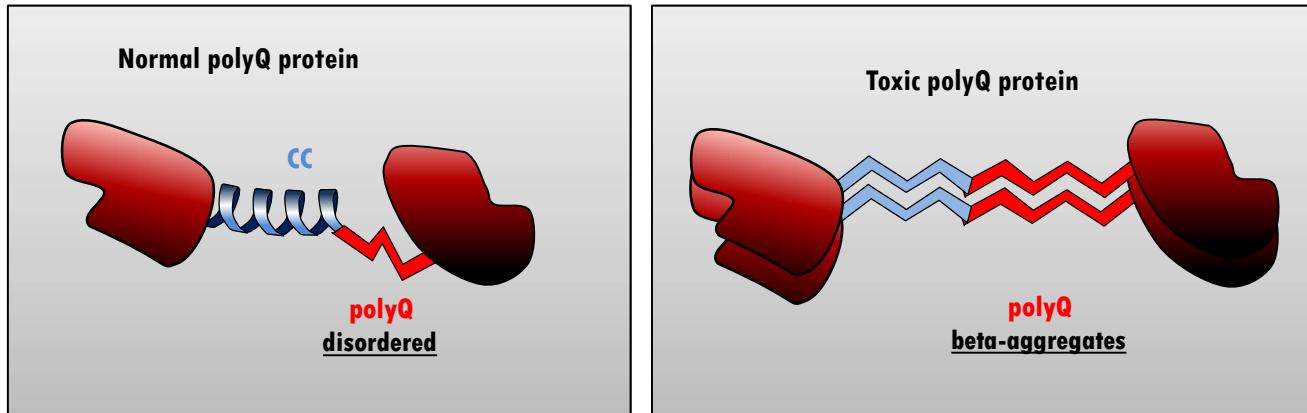


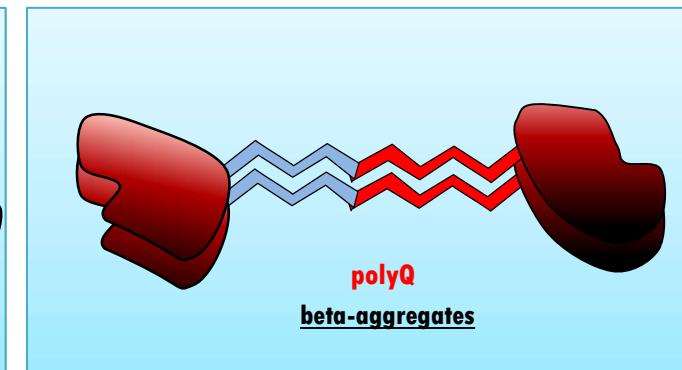
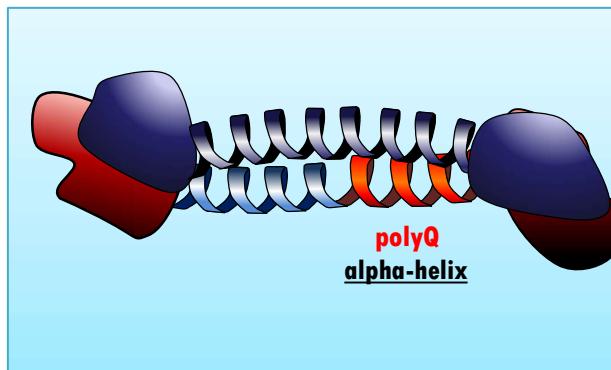
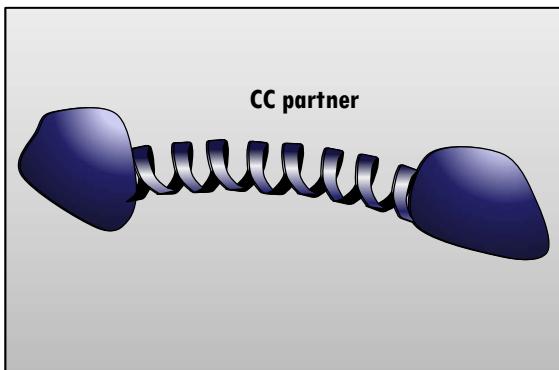
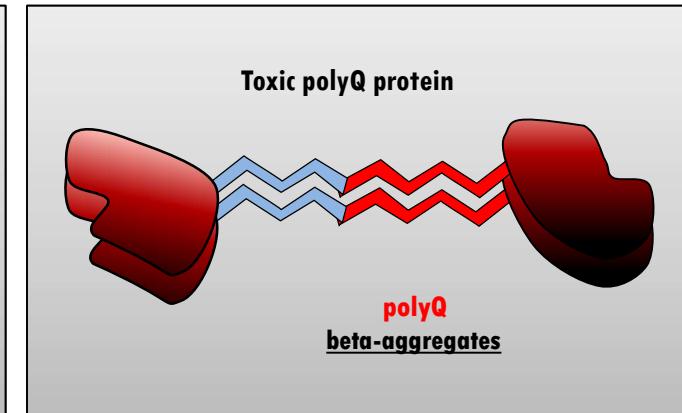
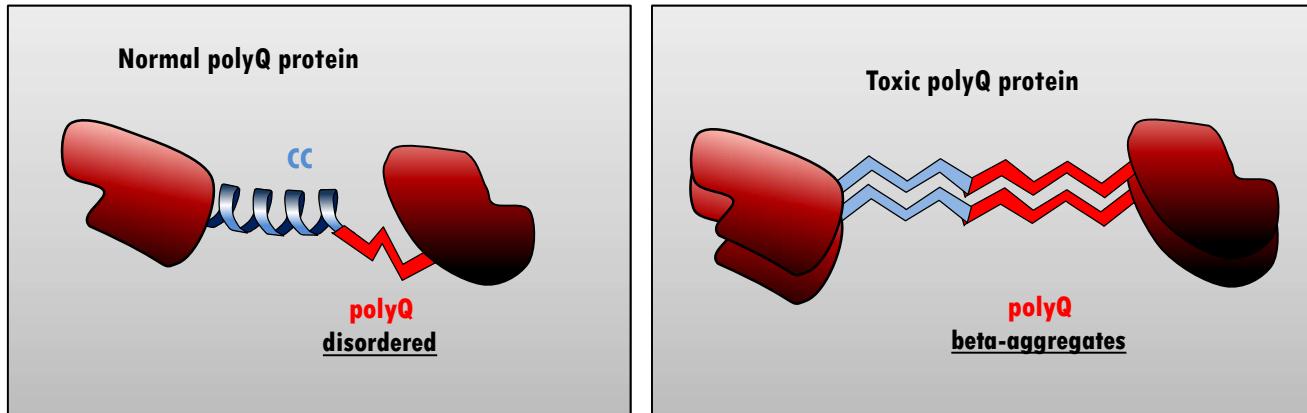


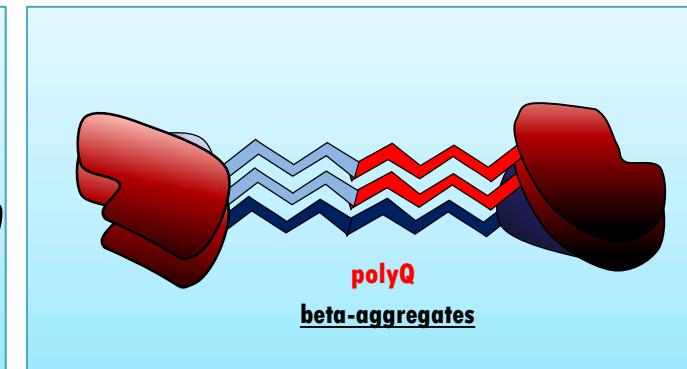
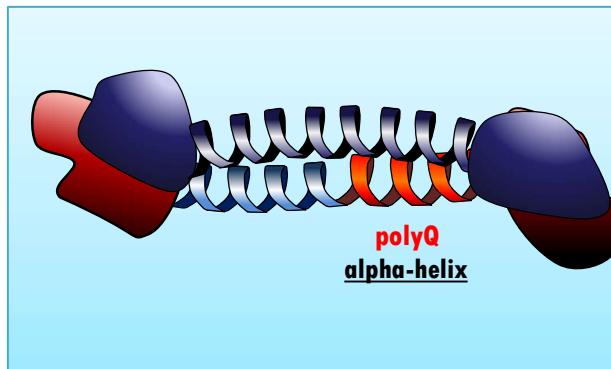
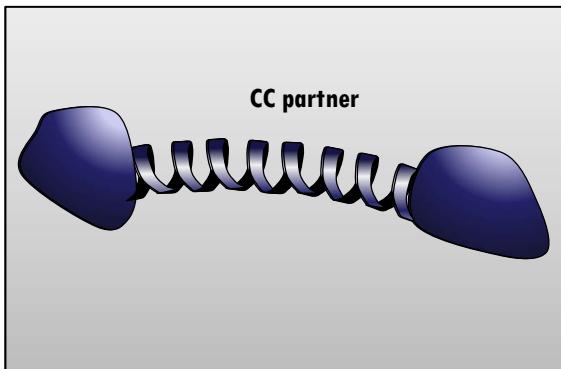
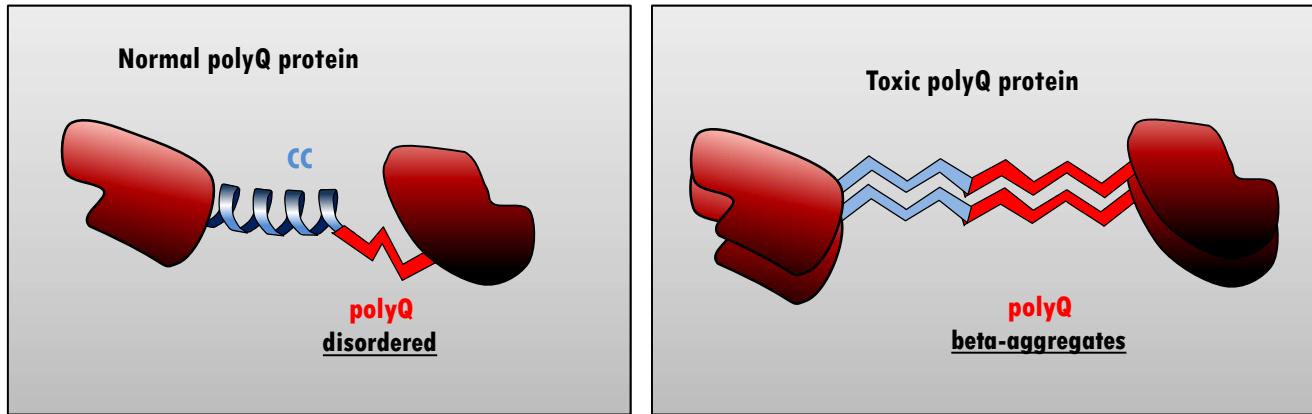


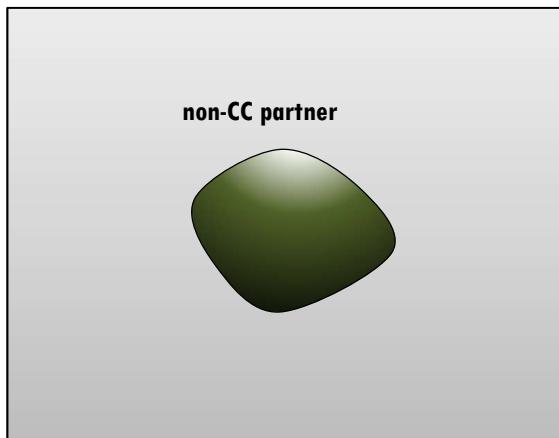
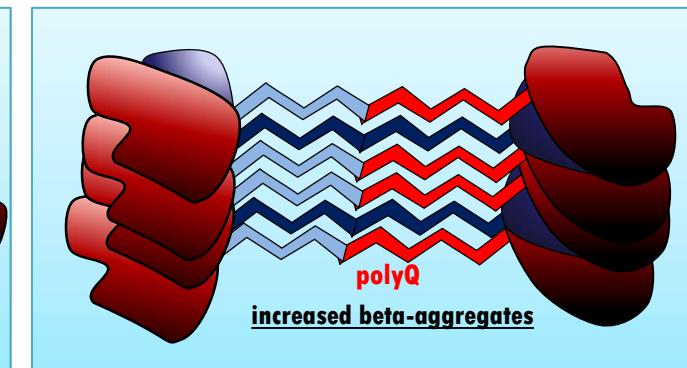
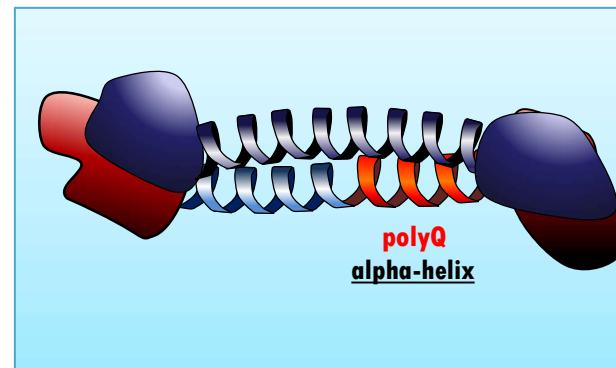
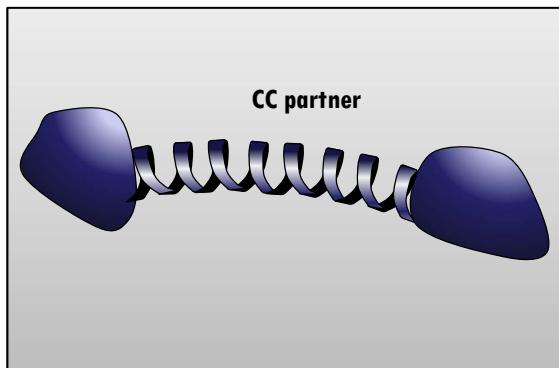
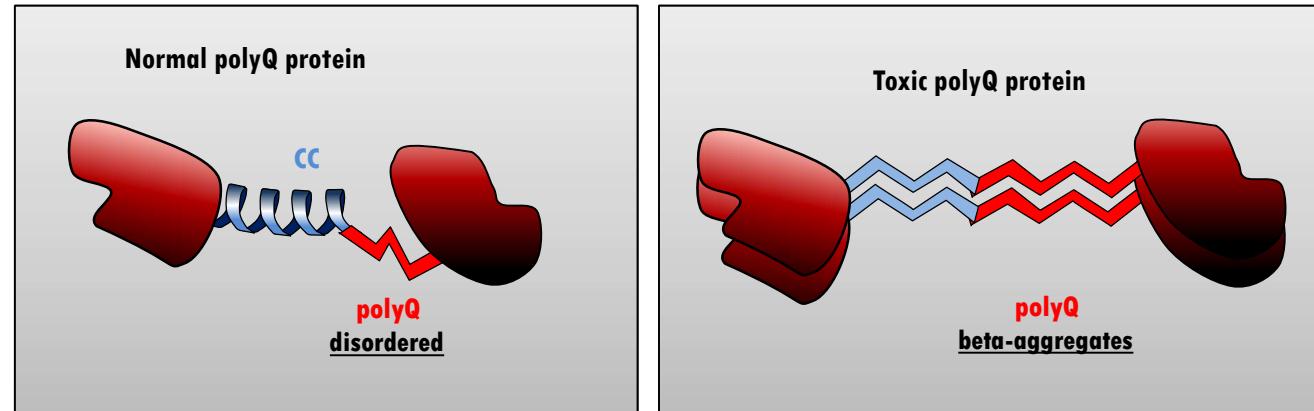


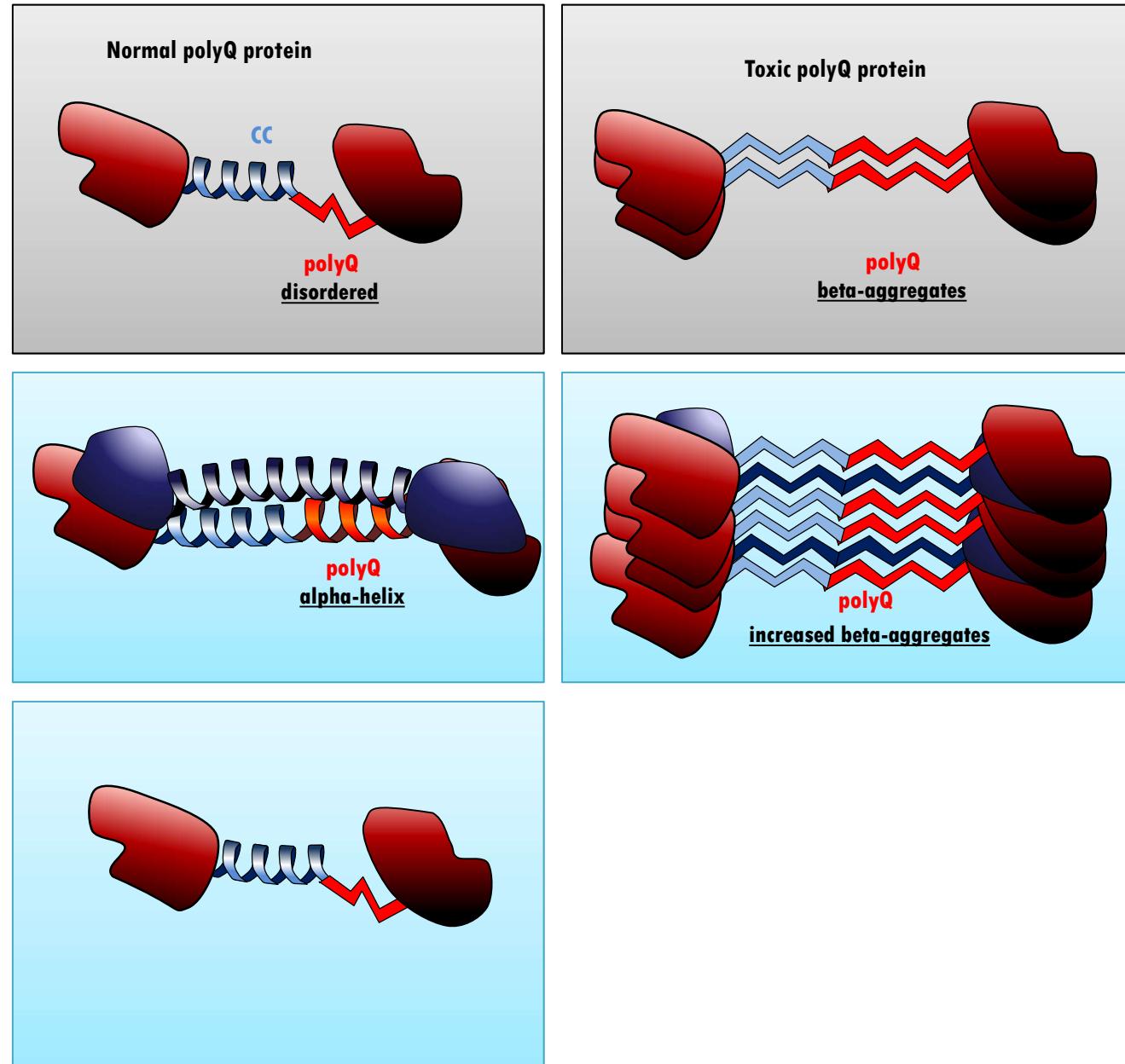


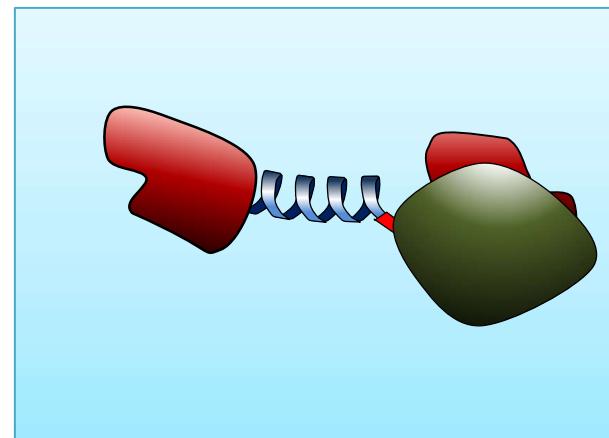
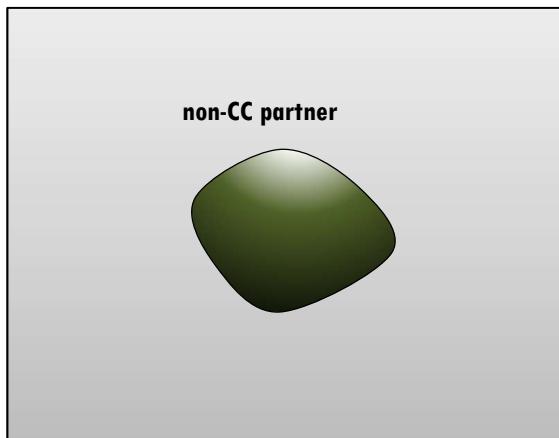
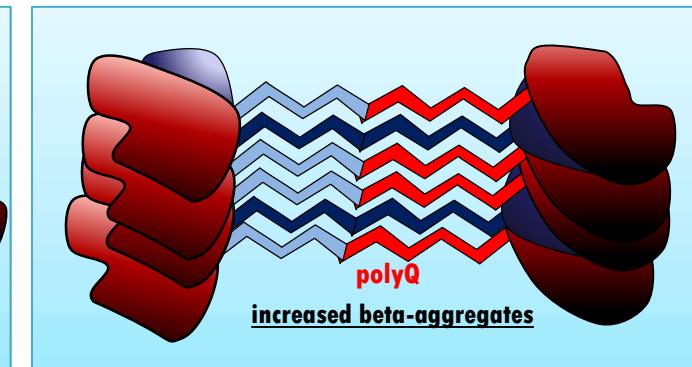
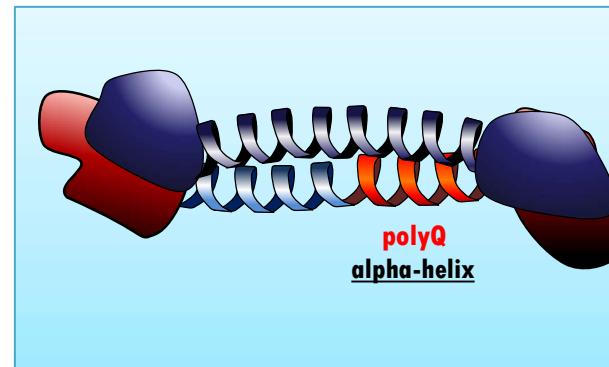
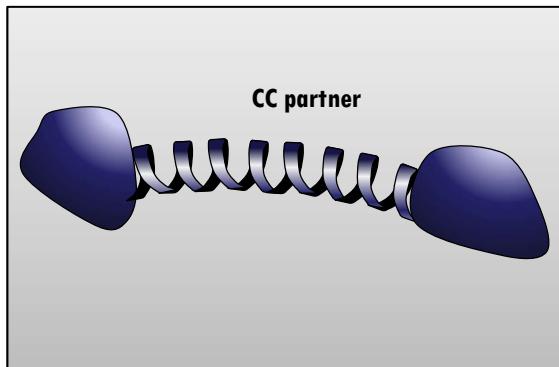
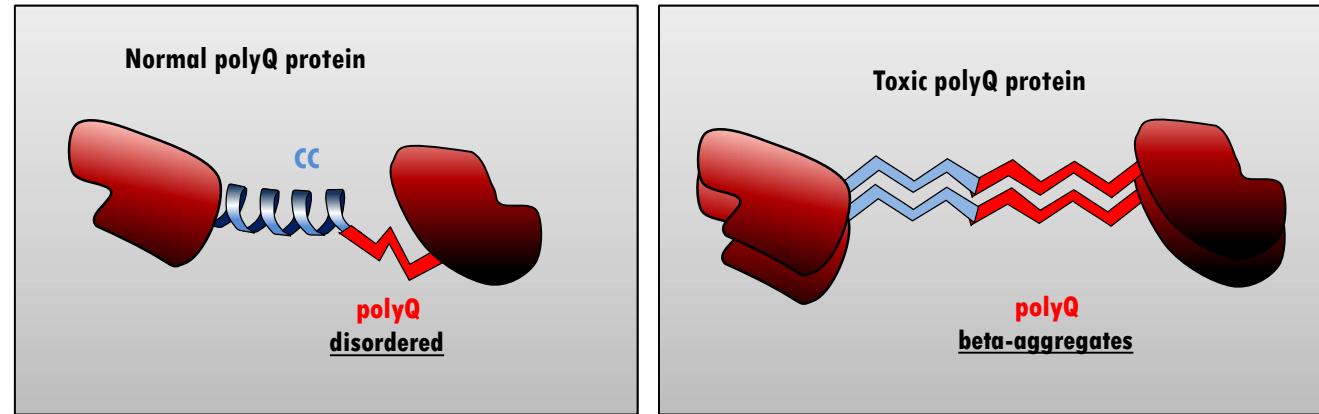


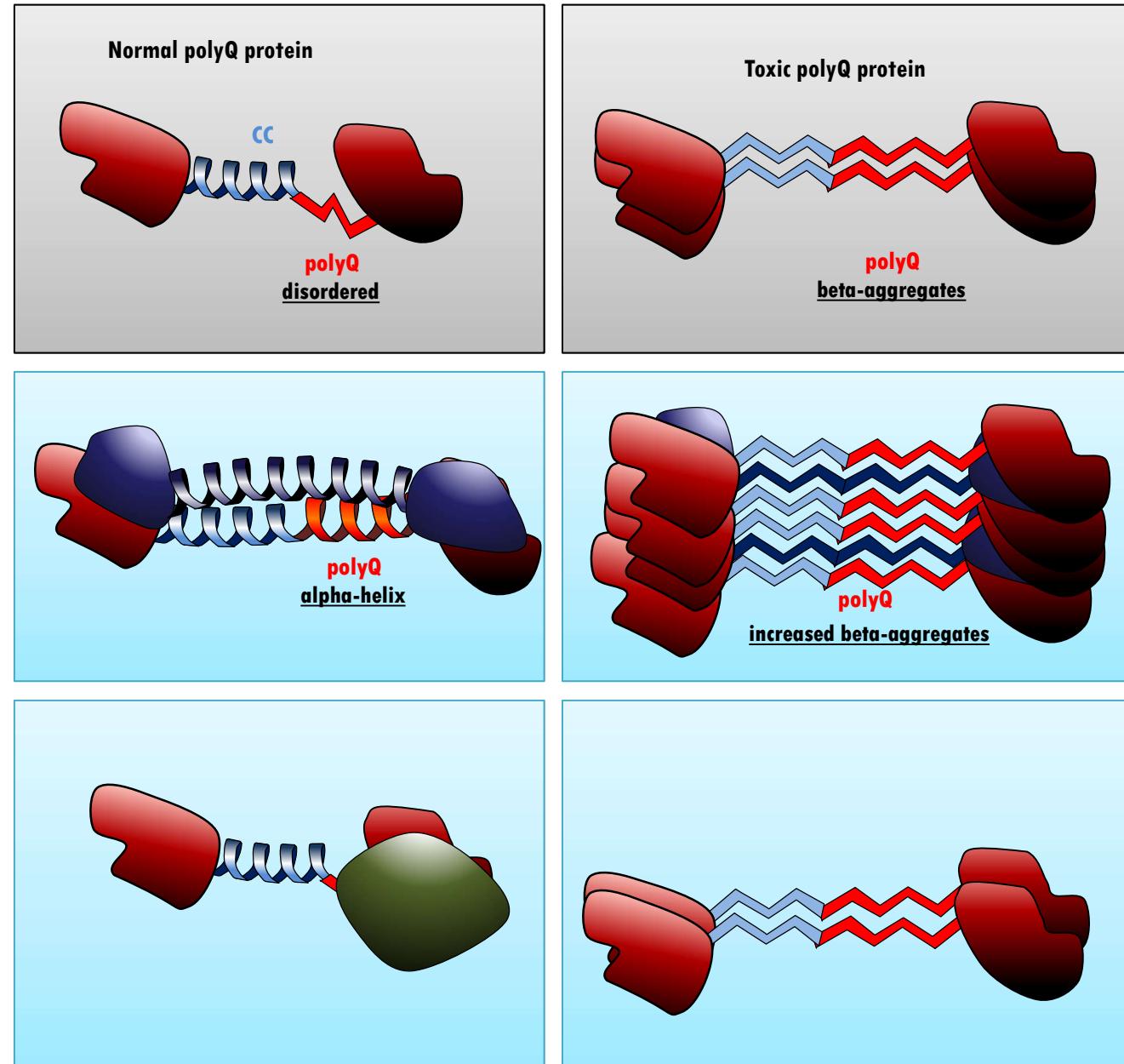


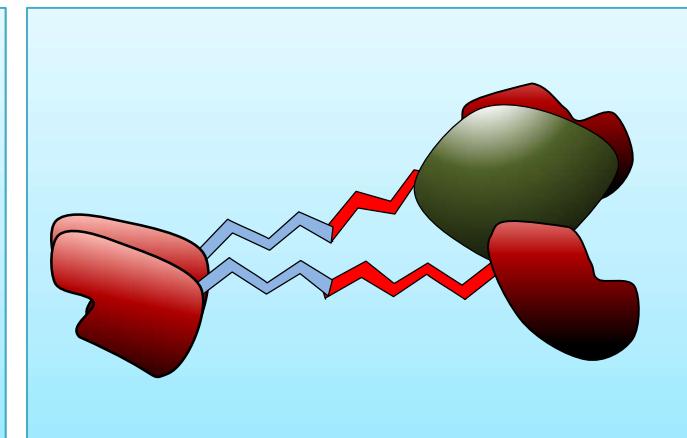
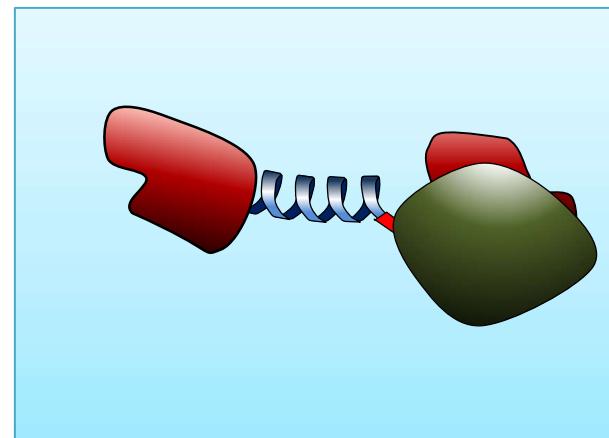
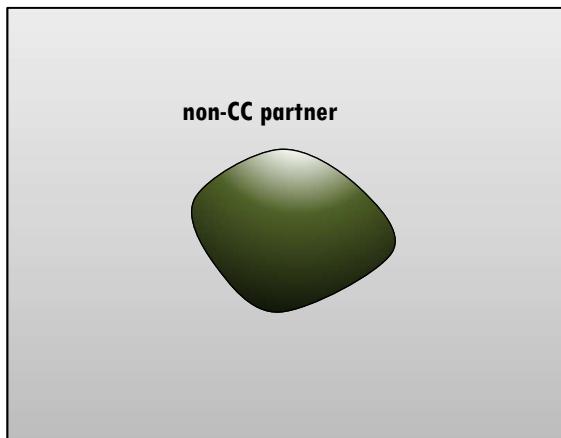
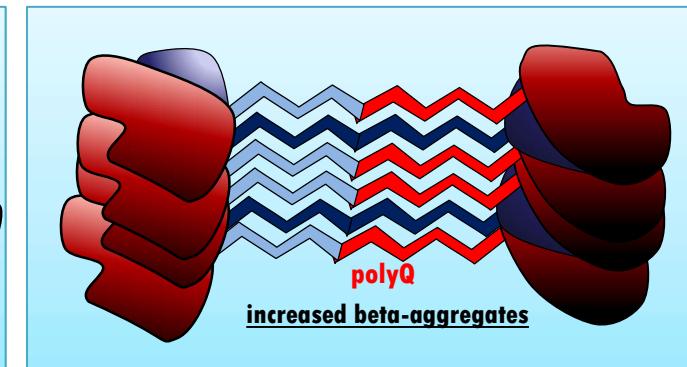
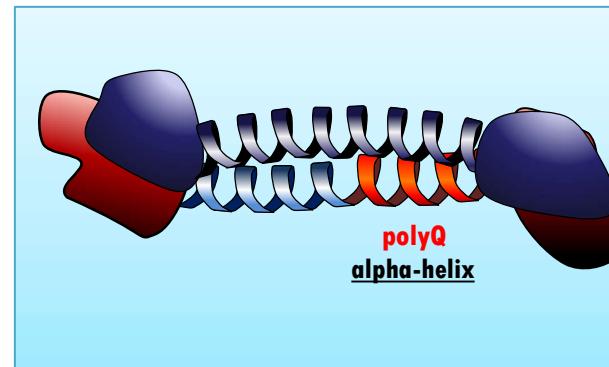
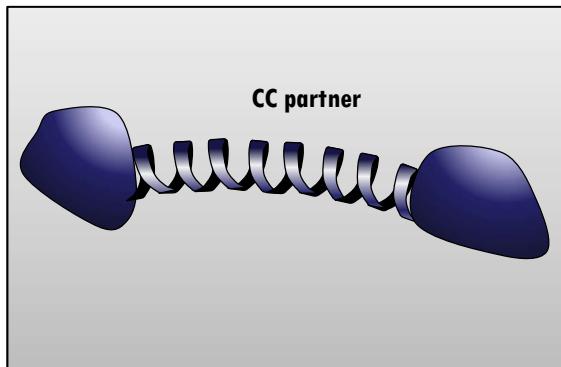
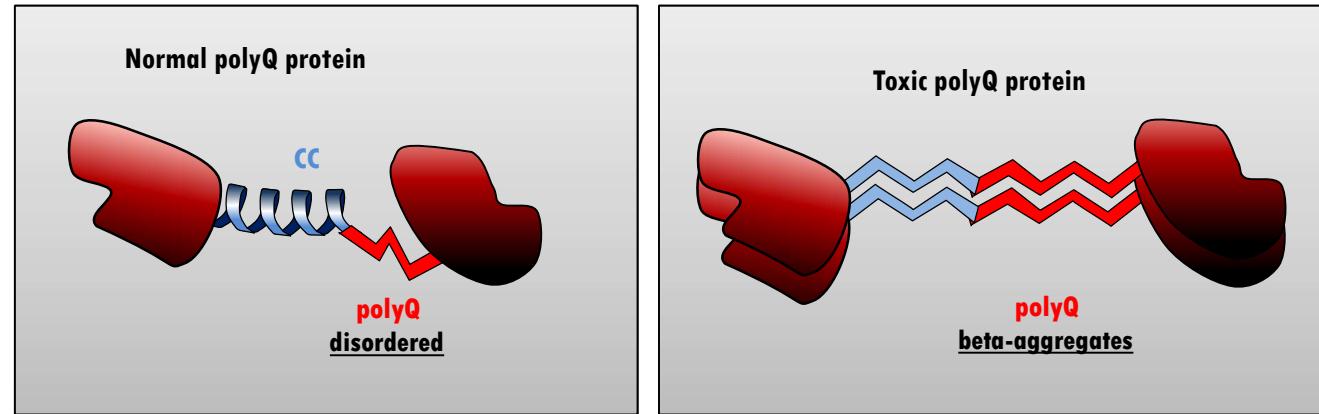


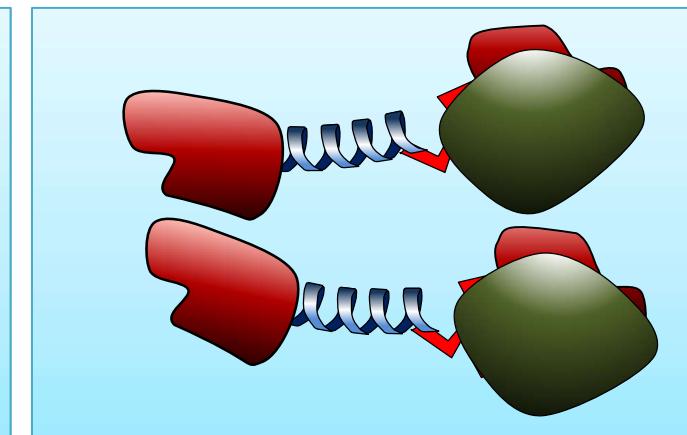
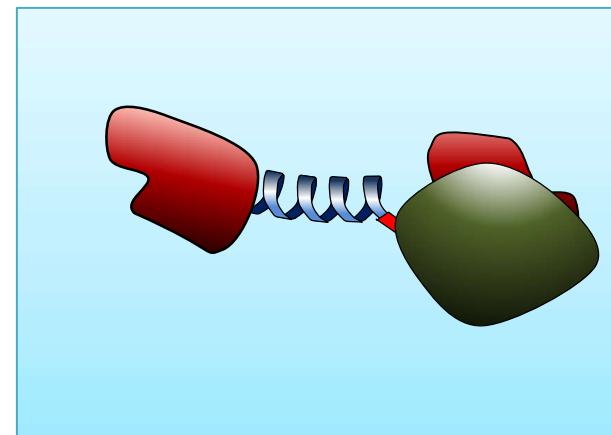
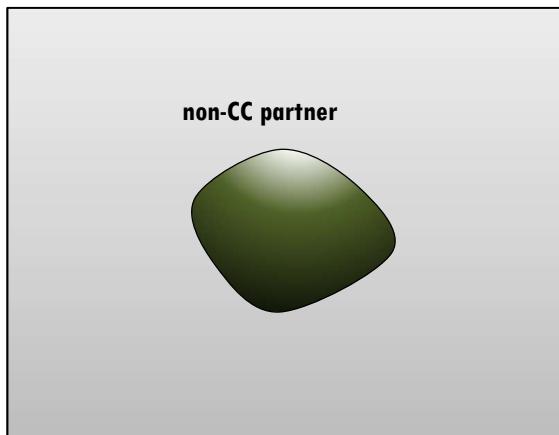
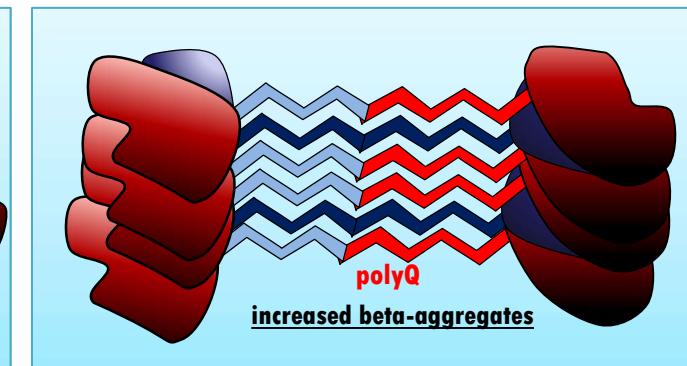
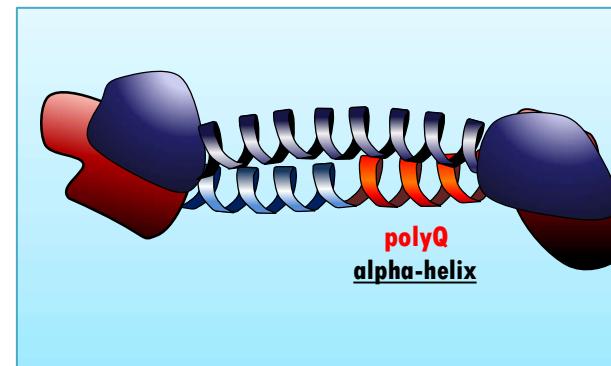
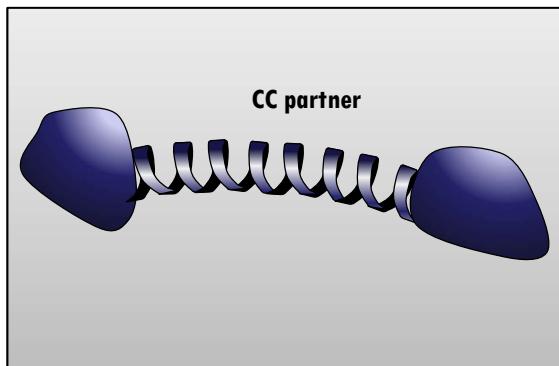
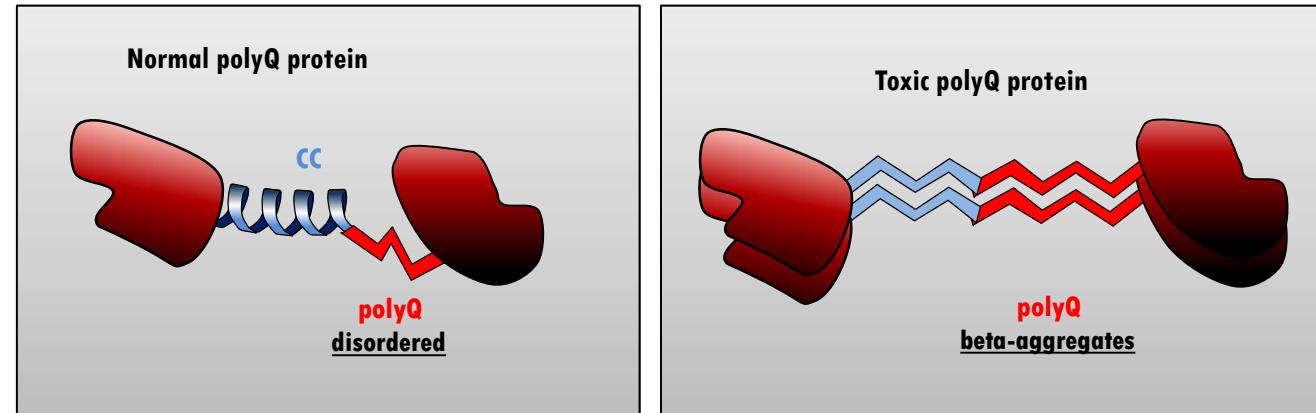












Exercise 3. Search for a polyQ insertion in the MR family

- Open in jalview the alignment of the mineralocorticoid receptor: MR1_fasta.txt
- Find a polyQ insertion.

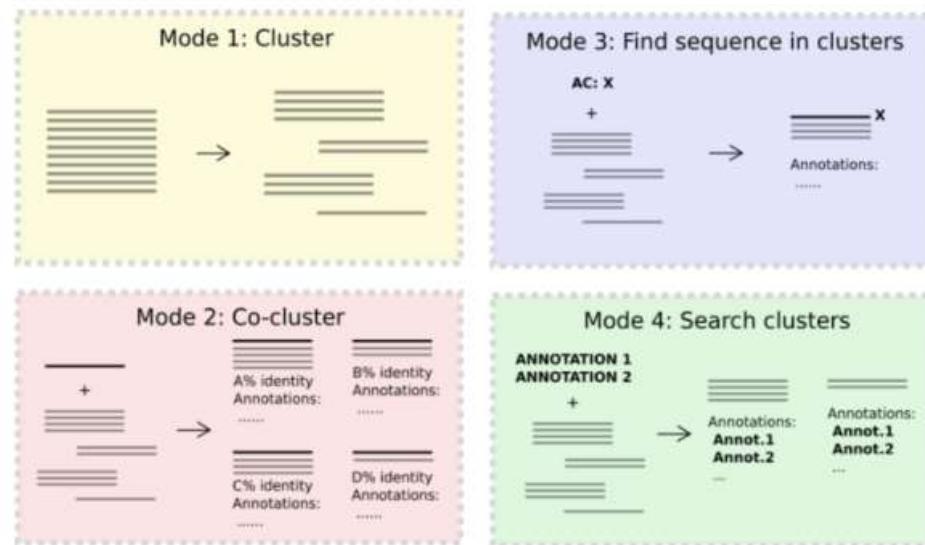
Do you see any other biased region nearby?

Clustering proteins

Pablo Mier



FastaHerder2



Mier and Andrade-Navarro (2016) *J. Comp. Biol.*

Clustering proteins

Pablo Mier



Results overview

Search settings

The cluster MUST have... The cluster MUST NOT have...

Organism/s = escherichia - PolyQ

Number of clusters found 2

Download link file .txt? 718580686328099.txt

Click on the leader to display its annotations

Leader: (1:)sp|A0K4S8|DNAK_BURCH
Leader: (1:)sp|Q83S00|FTSK_SHIFL

Time elapsed: 3 seconds

polyQ

...polyS regions? DM ...polyQ regions? DM
...polyG regions? DM ...polyA regions? DM
...polyL regions? DM ...polyM regions? DM
...polyW regions? DM ...polyY regions? DM

es separated by "+")

separated by "+")

In the cluster, ...

Escherichia

there **MUST** be at least one sequence from the following organism/s: (taxonomic id from an organism, e.g. 9606 for *H.sapiens*, or taxon name, e.g. *Homo*)*

there **MUST NOT** be any sequence from the following organism/s: (taxonomic id from an organism, e.g. 9606 for *H.sapiens*, or taxon name, e.g. *Homo*)*

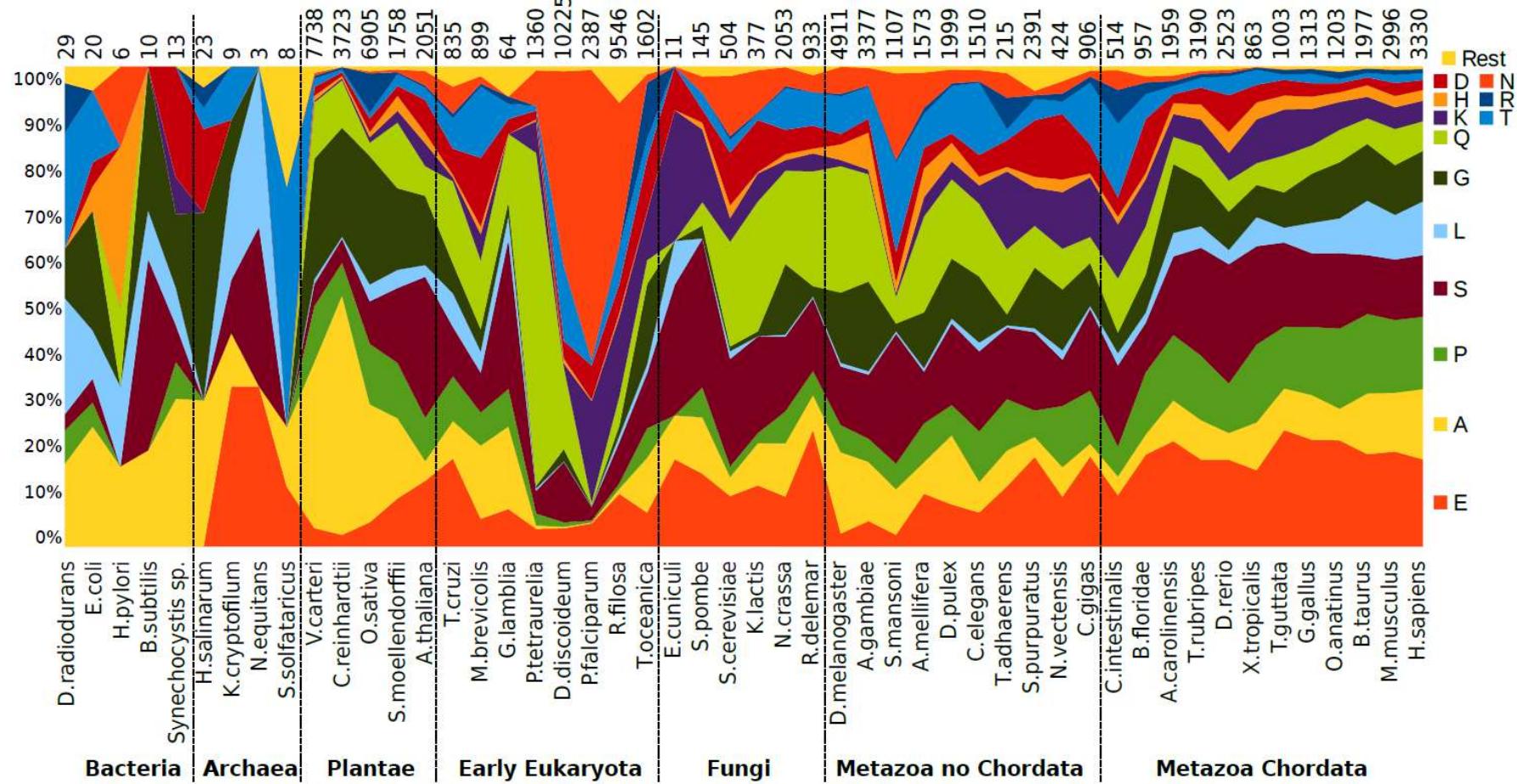
*if more than one, separate them by "+", e.g. 9606+Homo

SUBMIT

GO MODE 4! | What's this? | e.g. example 1, example 2

Mier and Andrade-Navarro (2016) *J. Comp. Biol.*

Frequency of homorepeats in 50 species



Mier et al. (2017) Proteins

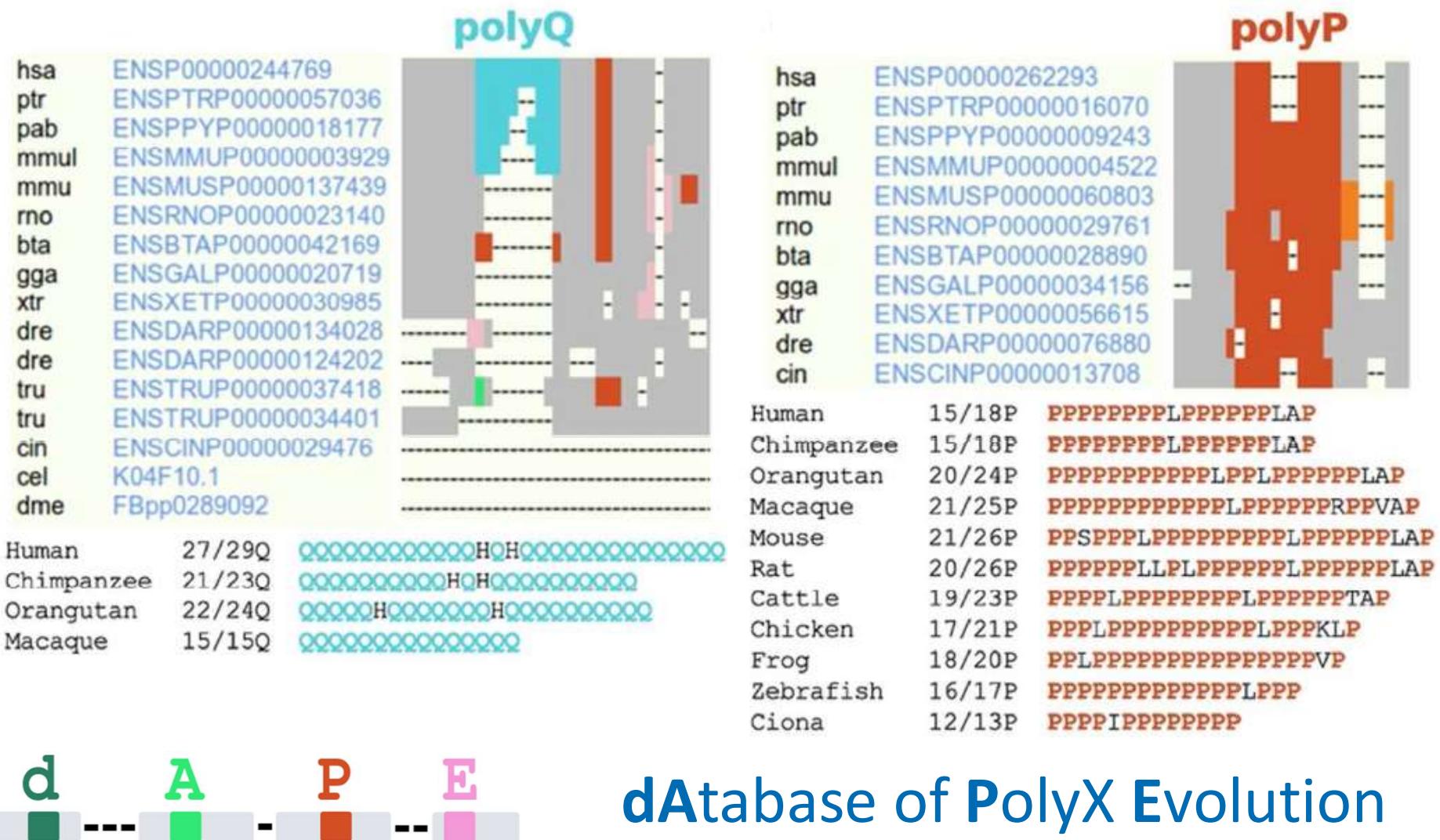
Context and evolution of homorepeats



dAtabase of PolyX Evolution

Mier *et al.* (2016) Bioinformatics

Context and evolution of homorepeats



Mier *et al.* (2016) Bioinformatics