

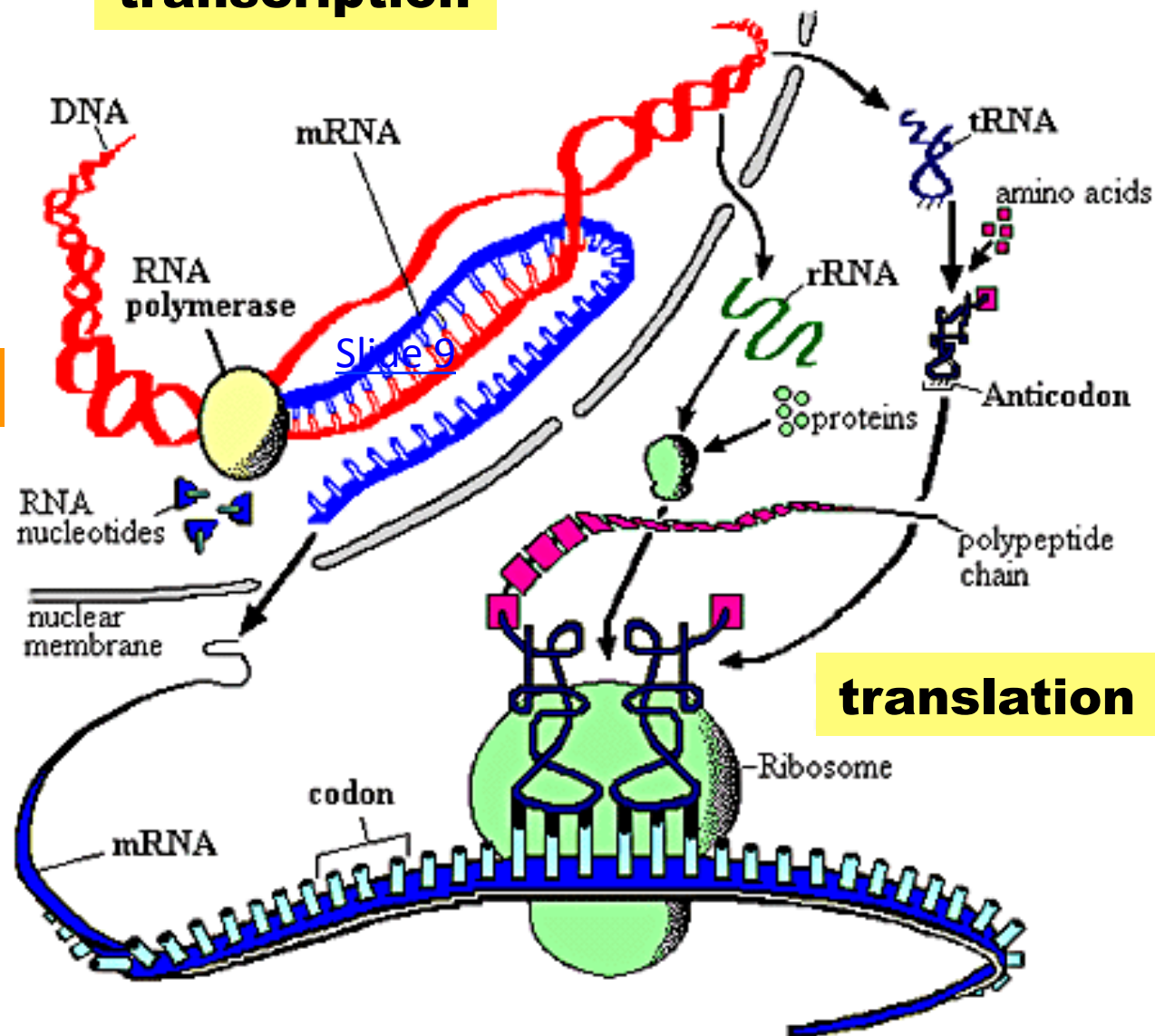
Assignment 1: Naïve Bayes Classifier of protein domains

Master in Bioinformatics UPF
2010-2011

Eduardo Eyras
Computational Genomics
Pompeu Fabra University - ICREA
Barcelona, Spain

transcription

nucleus

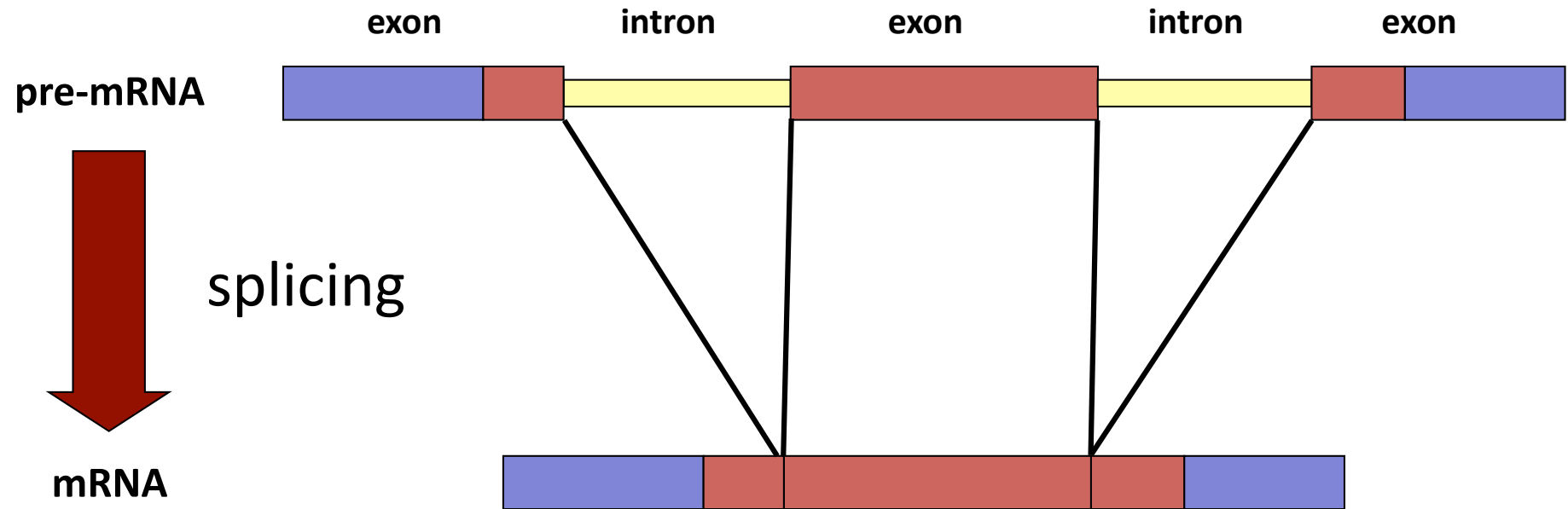


translation

cytoplasm

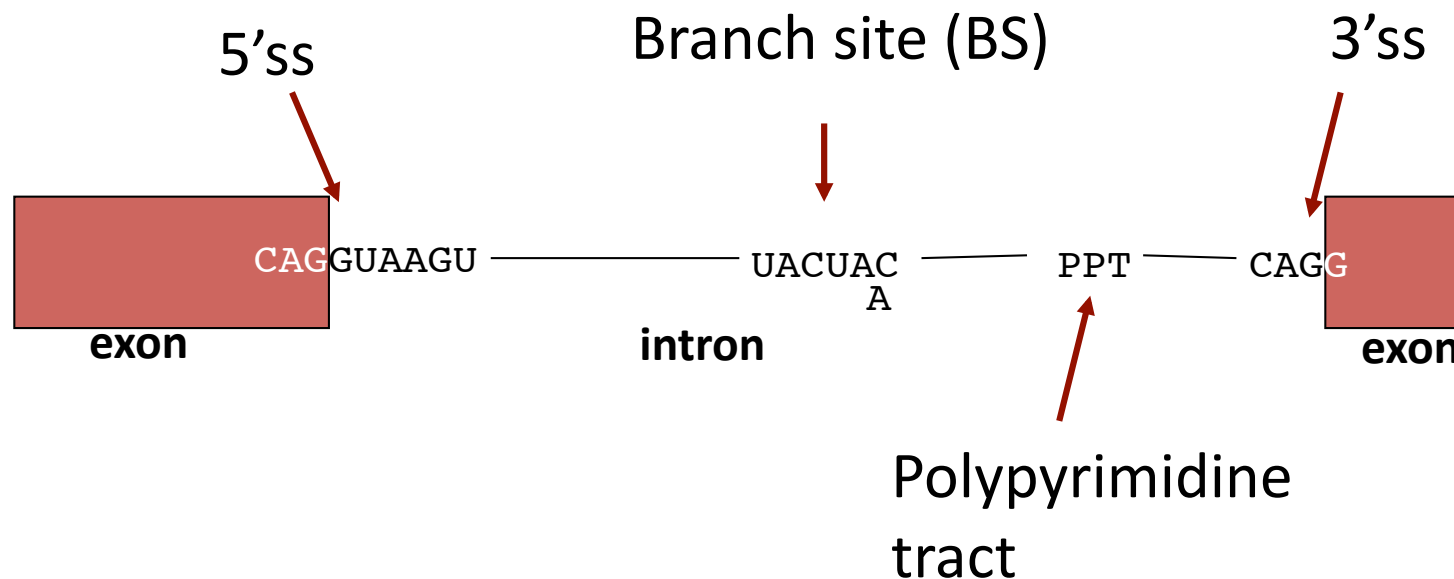
Protein synthesis

Splicing



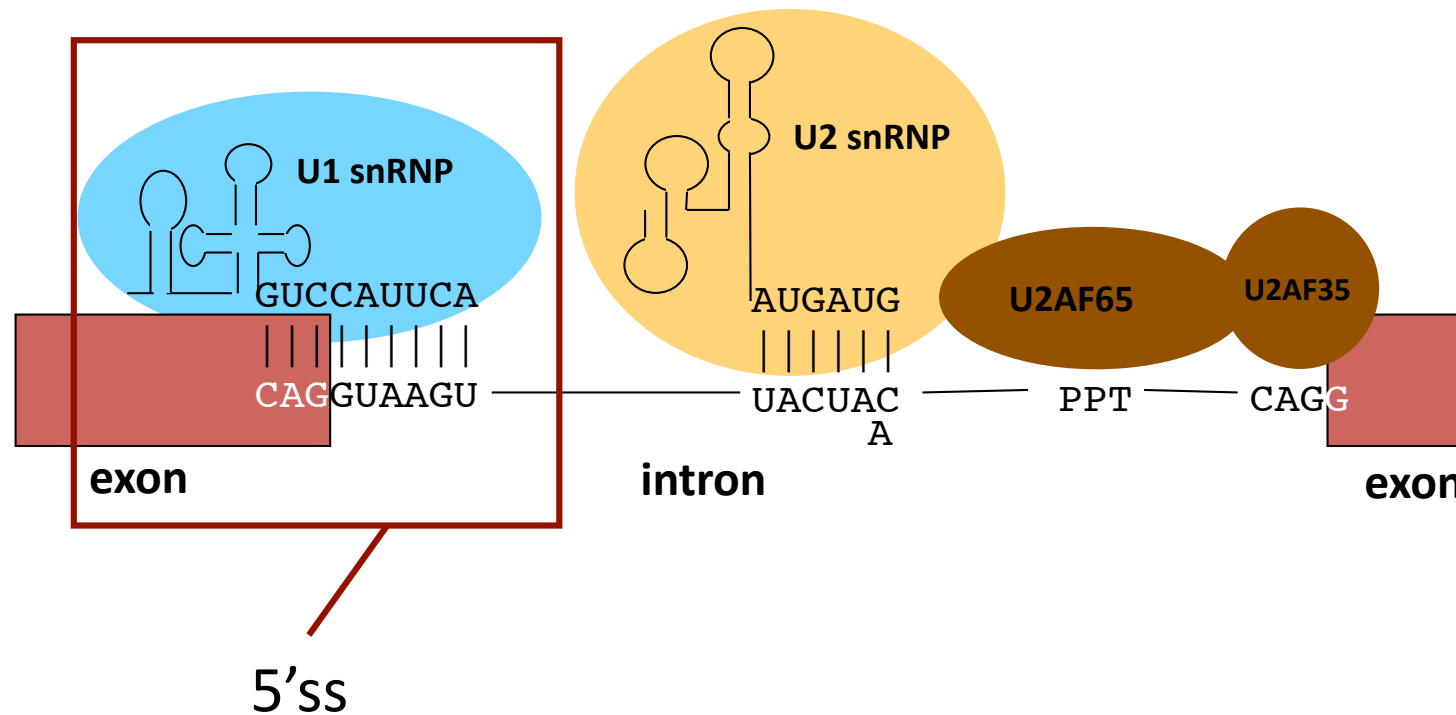
Introduction: splicing signals

4 basic splicing signals



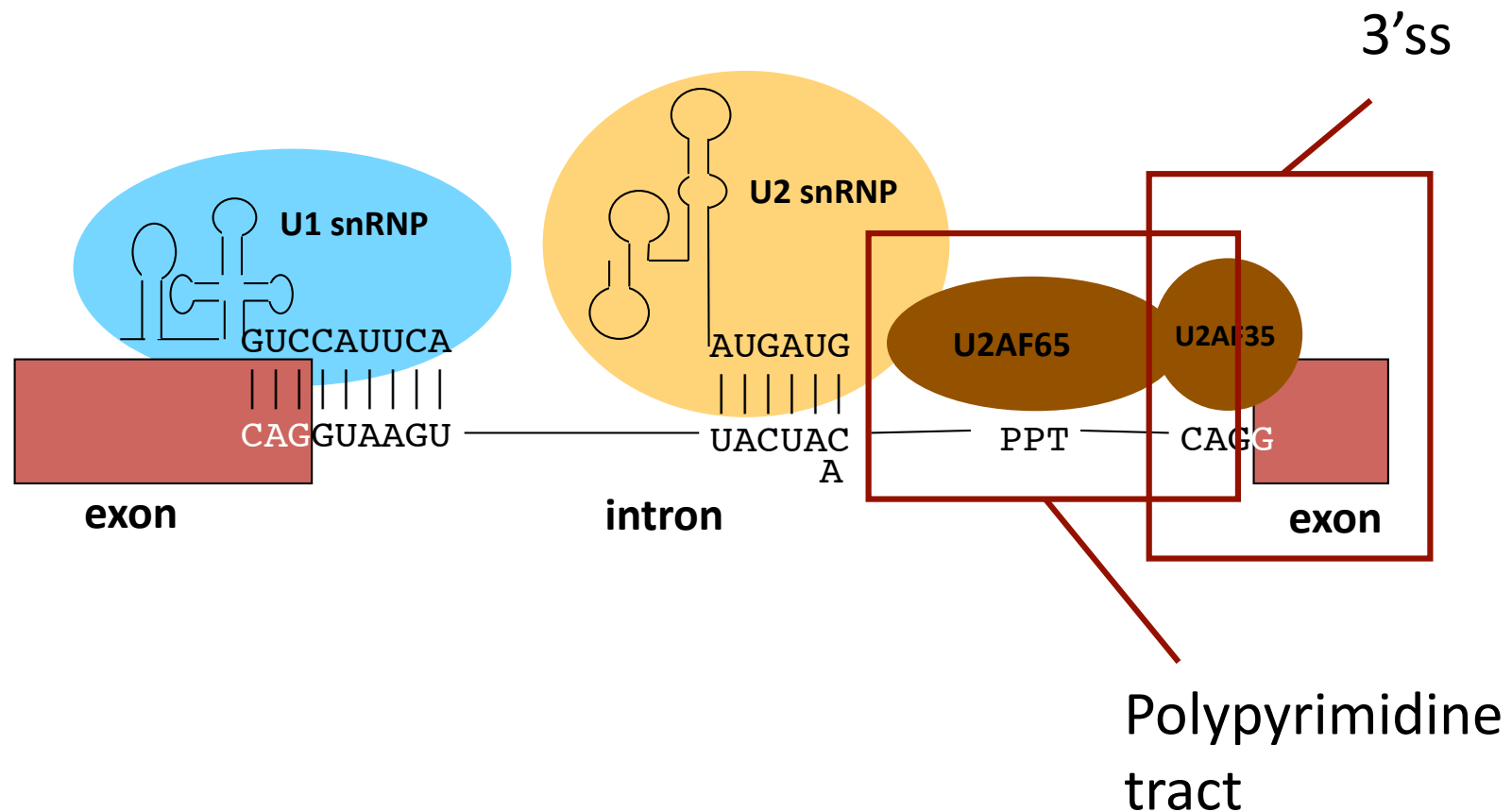
Introduction: splicing signals

4 basic splicing signals



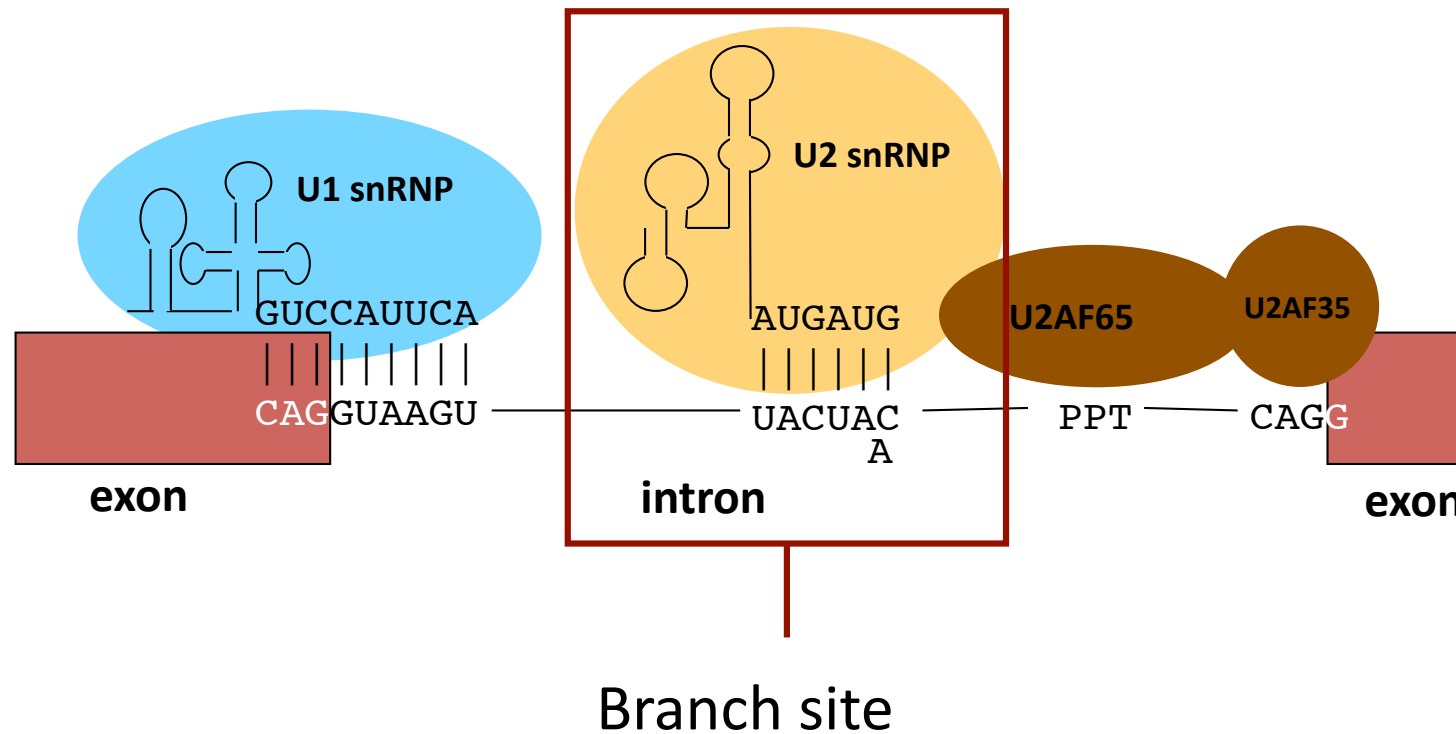
Introduction: splicing signals

4 basic splicing signals



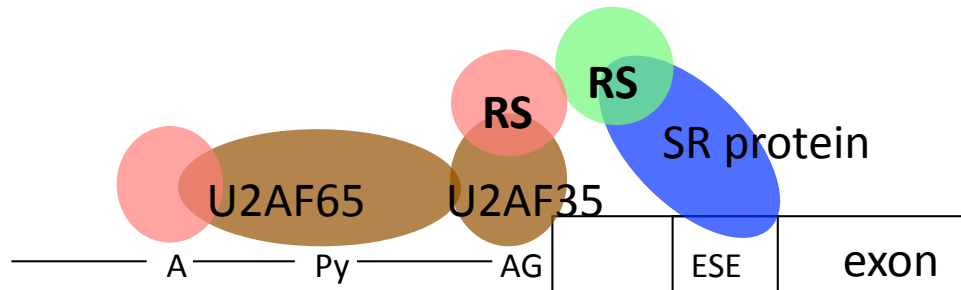
Introduction: splicing signals

4 basic splicing signals



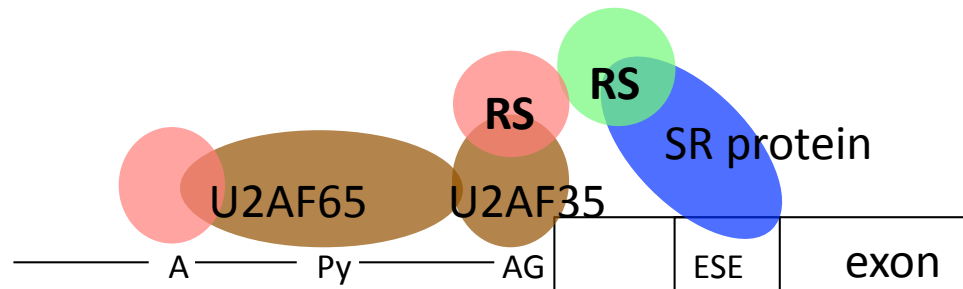
Splicing Factors (Protein-RNA interaction)

Splicing factors can bind to the pre-mRNA and promote (enhance) splice-site selection through the recruitment of spliceosomal components

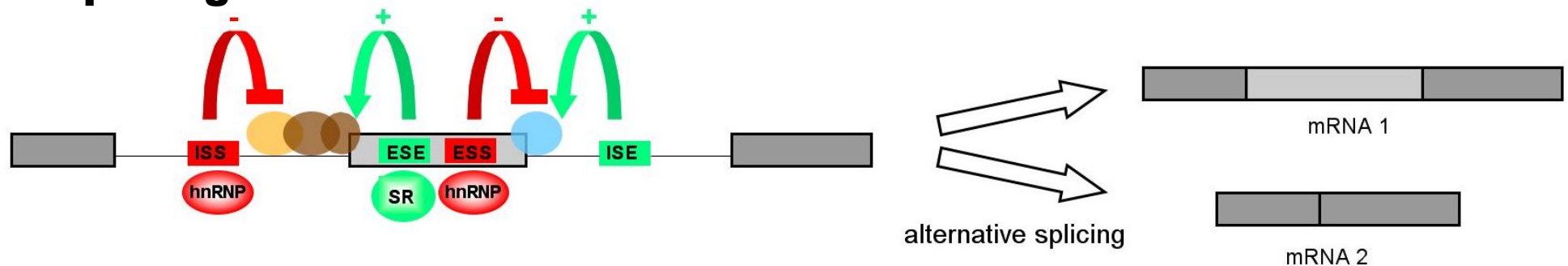


Splicing Factors (Protein-RNA interaction)

Splicing factors can bind to the pre-mRNA and promote (enhance) splice-site selection through the recruitment of spliceosomal components

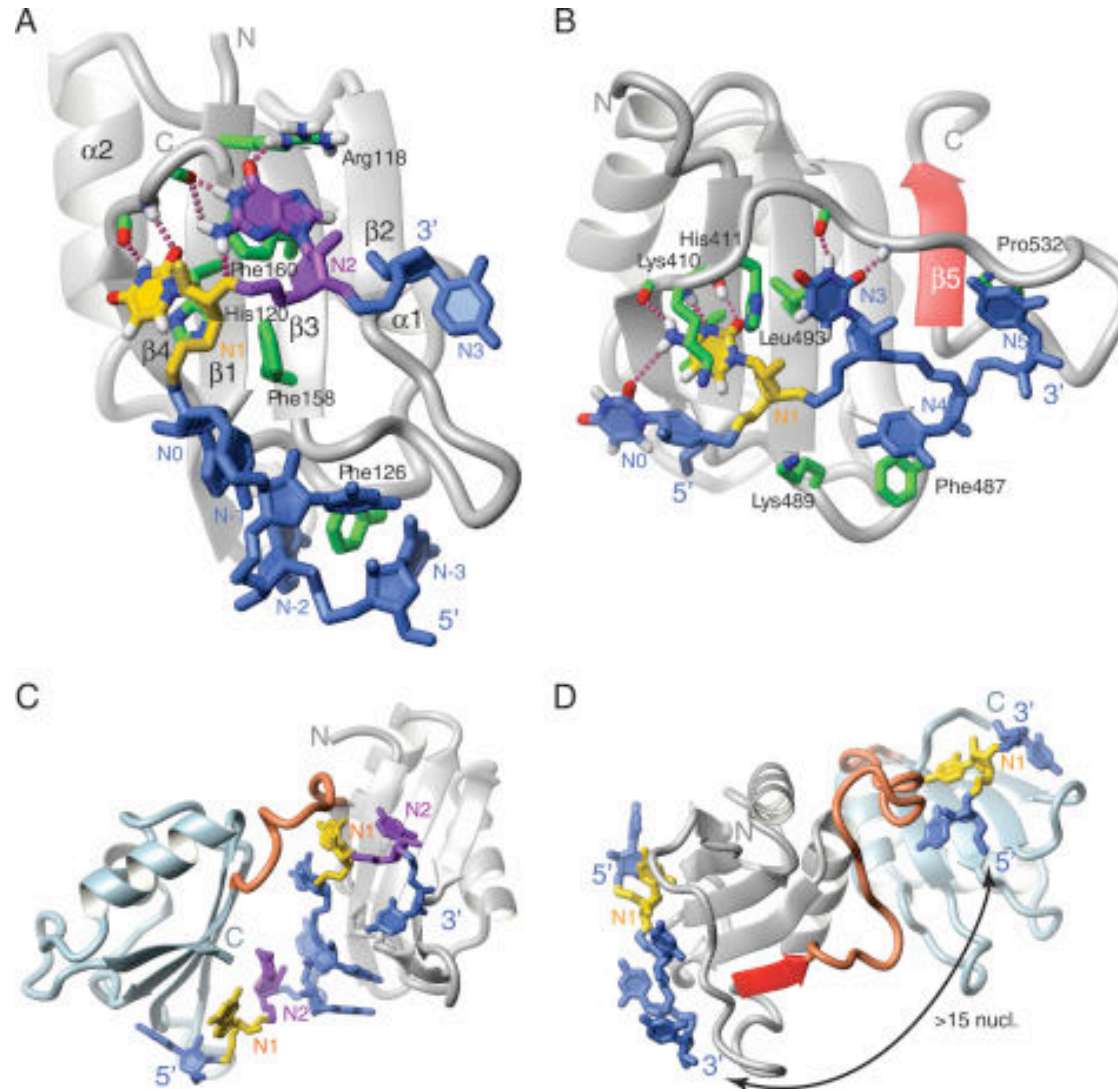


In fact, the regulated action of factors gives rise to alternative splicing



KH domains. (A) Type I KH domain of Nova (PDB code: 1EC6). (B) Type II KH domain of NusA (PDB code: 2ATW). (C) KH and QUA2 domains of SF1 (PDB code: 1K1G). (D) Tandem KH domains of NusA (2ATW)

RNA Binding Domains



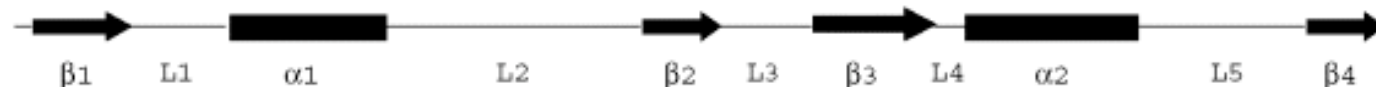
RRM domains. (A) The RRM of Fox-1 (PDB code: 2ERR). (B) RRM3 of PTB (PDB code: 2ADC). (C) The tandem RRM3 of Sex-lethal (PDB code: 1B7F). (D) RRM3 and 4 of PTB (PDB code: 2ADC)

The RNA Recognition Motif (RRM)

Hexamer RNP2 = [IVL] [FY] [IVL] X N L

Octamer RNP1 = [LR] G [FY] [GA] [FY] [VIL] X [FY]

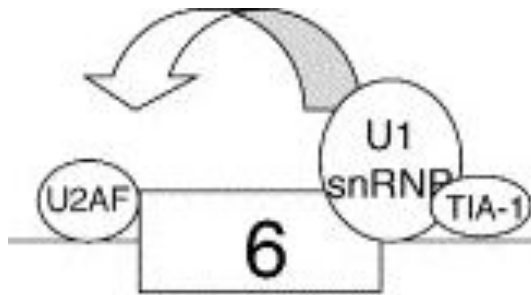
| | | | RNP2 | | | | | RNP1 | | | | | |
|-----------|--------|-----|---|--|---|--------------------------------|----------------|------|----|----|----|--|--|
| | | | | 10 | 20 | 30 | 40 | 50 | 60 | 70 | 80 | | |
| PTB | (1SJQ) | 60 | VIRKLPIDVTEGEVISLGLP---- | FGKVTNL----- | LMLKG----- | KNGAFIEMNTEEAANTMVNYT | SVTPVLRGQPIYIQ | 147 | | | | | |
| PTB | (1SRJ) | 183 | RIVVENLFPVTLDELH-QIFSK--- | FGTVLKI----- | ITFTKNN---- | QFQALLQYADPVSAQHAKLSLDG | QNIYNACCTLRID | 282 | | | | | |
| PTB | (1QM9) | 338 | VLLVSNLNPVTPQSLFILFGV---- | YGDVQRV----- | KILFNK----- | KENALVQADGNQAQLAMSHLNGHKLH-- | GKPIRIT | 407 | | | | | |
| PTB | (1QM9) | 455 | TLHLSNIPPSVSEEDLK-VLFSS--- | NGGVVKG----- | FKFPQKD---- | RKVALIQMGSVVEAVQALIDLHNHDLG- | ENHHLRVS | 531 | | | | | |
| Cstf-64 | (1P1T) | 17 | SVFVGNIPYEATEEQLK-DIFSE--- | VGPVVSF----- | RLVYDRETGKPKGYGFECEYQDQETALS | SAMRNLNGREPS-- | GRALRVD | 90 | | | | | |
| LA | (1OWX) | 244 | LKFSQDLDDQTCREDLHILFSNH--- | GEIK----- | WIDFVRGA-- | KEGIILFKEKAKEALGKAKDANNGNLQLRN | KEVTWEV | 305 | | | | | |
| TAP | (1FO1) | 121 | KITIPYGRKYDK-AWLLSMIQSKCSVPFTPIEFHYENTRAQFFVEDASTASAKAVNYKILDRENRRISIIINSSAP--- | PHS | 290 | | | | | | | | |
| ALY | (1N08) | 106 | KLLVSNLDFGVSDADIQ-ELFAE--- | FGTLKKA----- | AVHYDRSGR-SLGTADVHFERKADALKAMKQYNGVPLD-- | GRPMNIQ | 178 | | | | | | |
| hnRNP A1 | (1UP1) | 15 | KLFIGGLSFETTTDESLR-SHFEQ--- | WGTLTDC----- | VVMRDPNTKRSRGFGEVITYATVEEVDAAMNARP-HKVD-- | GRVVEPK | 87 | | | | | | |
| hnRNP A1 | (1HA1) | 105 | KIFVGGIKEDTEEHHLR-DYFEQ--- | YGKIEVI----- | EIMTDRGSGKKRGFAFVTPDDHDSVDKIVIQKY-HTVN-- | GHNCEVR | 177 | | | | | | |
| HUD | (1FXL) | 47 | NLIWNYLPQNMTOEEFR-SLFGS--- | IGEIESC----- | KLVRDKITGQSLGYGFVNYIDPKDAEKAINTLNGLRLQ-- | TKTIKV | 119 | | | | | | |
| HUD | (1FXL) | 133 | NLYVSGLPKMTMQELE-QLFSQ--- | YGRIITS----- | RILVDQVTGVSRGVGFIRFDKRIEABEAIKGLNQKPKSGATEPITVK | 206 | | | | | | | |
| SXL | (2SXL) | 126 | NLIWNYLPQDMTDRELY-ALFRA--- | IGPINTC----- | RIMRDYKTGYSYGYAFVDFVDFSEMSQRAIKVLNGITVR-- | NKRLKV | 199 | | | | | | |
| SXL | (1SXL) | 212 | NLYVTNLPRITITDDQLD-TIFGK--- | YGSIVQK----- | NILRDKLTGRPRGVAFVRYNKRREEAQEASALNNVIPEGGSQPLSVR | 290 | | | | | | | |
| PABP | (1CVJ) | 12 | SLYVGDLPDPVTEAMLY-EKFSP--- | AGPILSI----- | RVCRDMITRRSLGYVYVNFQQPADAEALDTPNFDVIK-- | GKPVRI | 84 | | | | | | |
| PABP | (1CVJ) | 99 | NIFIKNLDKSIDNKALYDTFSAF--- | GNILSCK----- | VVCDENGSKGYGFVHFETQEAERAIEKMNGMLLNDKRVFVGRPKS | 175 | | | | | | | |
| Nucleolin | (1FJE) | 309 | NLFIGNLNPNKSVAEKVAISEL--- | FAKND----- | LAVVDVRTGTNRKFGYVDFESAEDLEKAL-ELTGLKVF-- | GNEIKLE | 380 | | | | | | |
| Nucleolin | (1FJE) | 396 | LLAQNLNPNKSVAEKVAISEL--- | EIRLVSQ----- | DGKSKGIAYIEFKS-- | EADAEKNLEEKQGAID-- | GRSVSLY | 463 | | | | | |
| U1A | (1DZ5) | 11 | TIYINNLNEKIKKDELKKSLEYAI--- | FSQFGQI----- | LDILVSRSLKMRGQAFVIFKEVSSATNALRSMQGFPPY-- | DKPMRIQ | 85 | | | | | | |
| U2B" | (1A9N) | 8 | TIYINNMNDKIKKEELKRSLEYAL--- | FSQFGHV----- | VDIVALKTMKMRGQAFVIFKELGSSTNALRQLQGFPPY-- | GKPMRI | 81 | | | | | | |
| CBP20 | (1H2T) | 41 | TLVVGNLsfytteeqiy-ELFSK--- | SGDIKKI----- | INGLDKMKKTACGFCFVEYYSRADAENAMRYINGTRLD-- | DRIIRD | 114 | | | | | | |
| Y14 | (1P27) | 74 | ILFVTGVHHEATEEDIH-DKFAE--- | YGEIKNI----- | HLNLDRTTGYLKGYILVEYETYKEAQAAAMEGLNGQDLM-- | GQFISVD | 147 | | | | | | |
| UPF3 | (1UW4) | 52 | KVMIRRLPPTLTKEQLQEHLPQPM--- | PEHDYFE----- | FFSNDTSLYPHMYARAYINPKNQEDILFRDRPDGYVPLDNKGQEYPA | 131 | | | | | | | |
| U2AF65 | (1U2F) | 150 | RLVVGNIIPFGITEEAMM-DFFNAQMR-LGGLTQAPG--- | NPVLAVQINQDKMFLEPRSVDETTQAM-APDGIIPQ-- | GQSLKIR | 227 | | | | | | | |
| U2AF65 | (2U2F) | 260 | KLFIGGLPNYLNDQVK-ELLTS--- | FGPLKAF----- | NLVKDSATGLSKGYAFCEYVDINVTDAQIAGLNGMQLG-- | DKLLVQ | 333 | | | | | | |
| U2AF35 | (1JMT) | 66 | RSVSDVEMQEHYDEFFEEVFTMEEEKYGEVEEM---- | NVC-DNLGDHILVGMVYKFRREEDA | EKAVIDLNNRWFN-- | GQPIHA | 143 | | | | | | |



Aromatic residues for primary RNA binding: F, Y, (W,H)

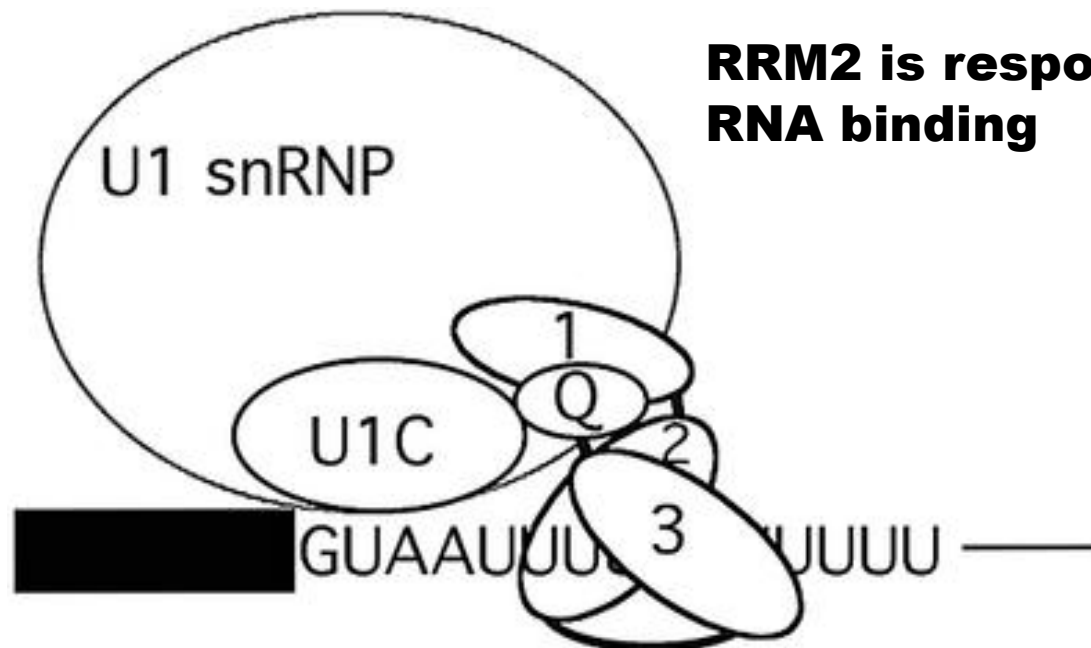
Other signatures: DAEDA (not found in all SR/SR-like proteins)

TIA1 can recruit U1 snRNP



Izquierdo et al. 2005

TIA1



RRM2 is responsible for RNA binding

Forch et al. 2002

TIA1 is conserved across eukaryotes

RRM1

(U1) interaction

| | 10 | 20 | 30 | 40 | 50 | 60 | 70 |
|------------------|-----------|-----------|------------|------------|----------|------------|-----------|
| CH99/1-80 | PHLYVGNLS | SPRVTDYIL | TEIFAVAGPV | VSAKIIQDRN | -FQHGGFN | YGFVEYAD | MRSADQAL |
| CJ21/1-80 | PHLYVGNLS | SPRVTDYIL | TEIFAVAGPV | VSAKIIQDRN | -FQHGGFN | YGFVEYAD | MRSADQAL |
| R265/1-80 | PHLYVGNLS | SPRVTDYIL | TEIFAVAGPV | VSAKIIQDRN | -FQHGGFN | YGFVEYAD | MRSADQAL |
| Aspergillus/1-81 | RALYVGG | LDPRTEDIL | KQIFETTGH | VISVKIIP | DKNQFNS | KGANYGFVE | FDDPGAAER |
| Rhizopus/1-81 | TTIYVGNL | DQRVTDML | NEIFTTVGG | VVSVKII | SVRKHNN | FGAVNYGFVE | FADPRVAEQ |
| Homo/1-77 | KTLYVGNL | SRDVTEALI | LQLFSQIG | PCKNCKMI | MDTA--- | GNDPYCFVE | FHEHRHAAA |

RRM2

(RNA binding)

| | 10 | 20 | 30 | 40 | 50 | 60 | 70 |
|------------------|---------|-----------|----------|-----------|---------|-----------|------------|
| CH99/1-79 | YHVFVGD | LSPEVNDDV | LSKAFGA | FGSLSEARV | MMDMNSG | KSRGYGFLS | FRDKADAEQA |
| CJ21/1-79 | YHVFVGD | LSPEVNDDV | LSKAFGA | FGSLSEARV | MMDMNSG | KSRGYGFLS | FRDKADAEQA |
| R265/1-79 | YHVFVGD | LSPEVNDDV | LSKAFGA | FGSLSEARV | MMDMNSG | KSRGYGFLS | FRDKADAEQA |
| Aspergillus/1-79 | FHIFVGD | LSNEVNDEV | LLQAFSA | FGSVSEARV | MMDMKTG | RSRGYGFVA | FRERADA |
| Rhizopus/1-79 | FHVFVGD | LAAEINDEK | LAQAFSE | FGTMSEAHV | MMDPLSG | KSRGFGFVA | FRDKTDA |
| Homo/1-79 | FHVFVGD | LSPQITTED | IKAAFAPE | GRISDARV | VKDMATG | KSKGYGFVS | EFNKVDA |

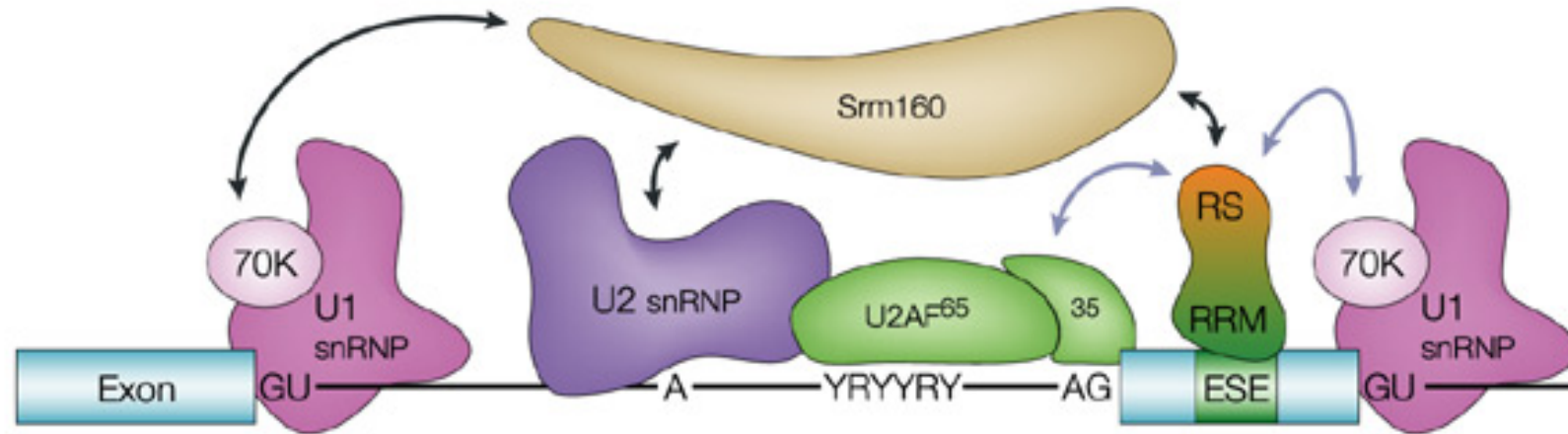
RRM3

| | 10 | 20 | 30 | 40 | 50 | 60 | 70 |
|------------------|-----------|-----------|------------|----------|----------|-----------|-------------|
| CH99/1-73 | TTVYVGNL | IPIYTTQAD | LIPLFQGYGY | IVEIRMQA | DRGFAFVK | LDTHQNAAL | AI THLQNLVH |
| CJ21/1-73 | TTVYVGNL | IPIYTTQAD | LIPLFQGYGY | IVEIRMQA | DRGFAFVK | LDTHQNAAL | AI THLQNLVH |
| R265/1-73 | TTVYVGNL | IPIYTTQAD | LIPLFQGYGY | IVEIRMQA | DRGFAFVK | LDTHQNAAL | AI THLQNLVH |
| Aspergillus/1-73 | TTTCYVGNL | TPYTTQND | IVPLFQNF | GYVIE | ETRMQA | DRGFAFIK | MDTHENAASA |
| Rhizopus/1-73 | TSIYVGNL | IPLNVSQND | LVQPFQR | FGYVQEV | KFQADR | GFAFVKMD | THENAANA |
| Homo/1-73 | CTVYCGVT | SGLTEQLMR | QTFSPFGQ | IMEIRVFP | DKGYSFVR | FNSHESAA | HAIVSVNGT |

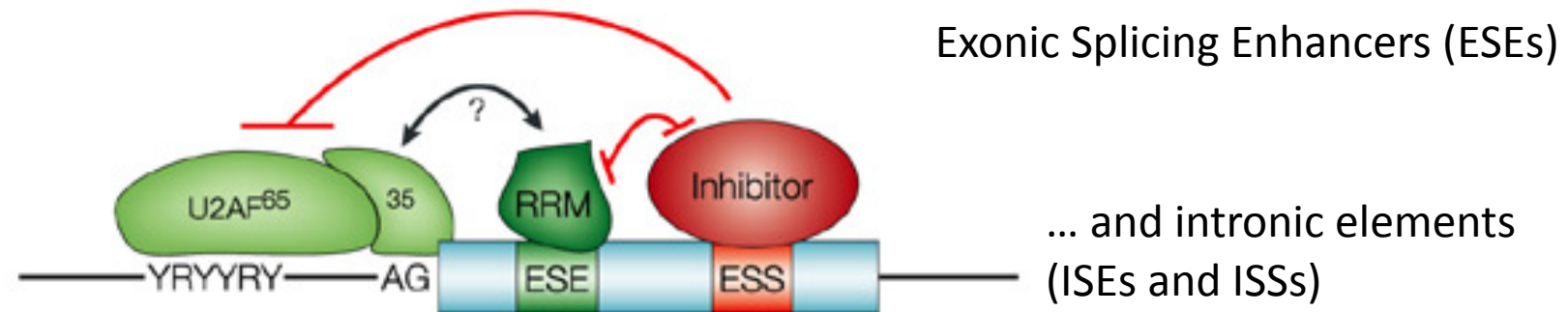
Functionally related proteins have also conserved RRM domains

| | | RNP2 | | | | | | | | | | RNP1 | | | | | | | | | |
|------|----------------------------|-----------|---------|---------|-----------|---------|--------|-------|---------|---------|--------|--------|--------|--------|--------|--------|--------|-------|-----|-------|-----|
| | | | 10 | | 20 | | 30 | | 40 | | 50 | | 60 | | 70 | | | | | | |
| TIA1 | <i>H_sapiens_TIA1</i> | FHVFVGDLS | PEITTED | EDIKAAF | -APFG-RIS | SDARVVK | DMA | TGKS | KGYGFVS | FNK | WDAENA | IQQMG | QWLGG | RQ | IRTNW | ATR | | | | | |
| | <i>M_musculus_TIA1</i> | FHVFVGDLS | PEITTED | EDIKAAF | -APFG-RIS | SDARVVK | DMA | TGKS | KGYGFVS | FNK | WDAENA | IQQMG | QWLGG | RQ | IRTNW | ATR | | | | | |
| | <i>C_familiaris_TIA1</i> | FHVFVGDLS | PEITTED | EDIKAAF | -APFG-RIS | SDARVVK | DMA | TGKS | KGYGFVS | FNK | WDAENA | IQQMG | QWLGG | RQ | IRTNW | ATR | | | | | |
| | <i>G_gallus_TIA1</i> | FHVFVGDLS | PEITTED | EDIKAAF | -APFG-RIS | SDARVVK | DMA | TGKS | KGYGFVS | FNK | WDAENA | IQQMG | QWLGG | RQ | IRTNW | ATR | | | | | |
| | <i>D_rerio_TIA1</i> | FHVFVGDLS | PEISTDD | VRAAF | -APFG-KIS | SDARVVK | DLAT | TGKS | KGYGFIS | FNK | WDAESA | IQQMNG | QWLGG | RQ | IRTNW | ATR | | | | | |
| | <i>D_melanogaster_ROX8</i> | HHIFVGDLS | PEIETET | LREAF | -APFG-EIS | NCRIVR | DPT | MTKS | KGYAFVS | FNK | WDAENA | IQAMNG | QWIGS | RS | IRTNW | STR | | | | | |
| | <i>C_elegans_TIA1</i> | FHVFVGDLS | SEVDNQK | LREAF | -QPF | FG-DVSD | AKVIR | DTNT | TKS | KGYGFVS | FNK | WDAENA | IQQMNG | QWLGG | RRT | IRTNW | ATR | | | | |
| PUB1 | <i>C_neoformans_PUB1</i> | YHVFVGDLS | PEVND | DVLSKAF | -GAF | GSLS | EARVMW | DMNS | GKSR | GYGFL | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | | |
| | <i>U_maydis_PUB1</i> | SHLFVGDLS | PDVDD | ALLQSSF | -SRFT | SLADVR | VMYDA | ATGK | SRGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | | | |
| | <i>A_fumigatus_PUB1</i> | FHIFVGDLS | NEVNDEI | LLQAF | -SAFG | SVSEAR | VMWDM | KTGR | SRGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | | | |
| | <i>M_grisea_PUB1</i> | FHIFVGDLS | NEVNDEI | LLQAF | -SAFG | SVSEAR | VMWDM | KTGR | SRGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | | | |
| | <i>N_crassa_PUB1</i> | FHIFVGDLS | NEVNDEI | LLQAF | -SAFG | SVSEAR | VMWDM | KTGR | SRGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | | | |
| | <i>Y_lipolytica_PUB1</i> | YNLFVGDLS | GADVND | EMLHKHF | -AHIP | GLLDAR | VMWDM | TTGR | SRGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | | | |
| | <i>D_hansenii_PUB1</i> | FNIFVGDLS | PEVDD | ETLNKSF | -SKFP | SLKQAH | VMWDM | QTGR | SRGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | | | |
| | <i>A_gossypii_PUB1</i> | FNLFVGDLS | NVDVDD | ETLSSTF | -KEFP | TFIQA | HVMWDM | QSGR | SRGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | | | |
| | <i>K_lactis_PUB1</i> | FNLFVGDLS | NVDVDD | ETLAGTF | -KEFP | TFIQA | HVMWDM | QSGR | SRGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | | | |
| | <i>C_glabrata_PUB1</i> | FNLFVGDLS | NVDVDD | ETLAGTF | -REFP | TFIQA | HVMWDM | QSGR | SRGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | | | |
| | <i>S_cerevisiae_PUB1</i> | FNLFVGDLS | NVNDD | ETLRNAF | -KDFP | SYLSGH | VMWDM | QTGR | SRGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | | | |
| | <i>A_thaliana_OBa</i> | FNIFVGDLS | PEVTD | AMLF | FTCF | -SVYP | TCS | DA | VMWDM | QKTGR | SRGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | |
| | <i>A_thaliana_OBb</i> | FNIFVGDLS | PEVTD | AA | LFDSF | -SAFN | SCS | DA | VMWDM | QKTGR | SRGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | |
| | <i>A_thaliana_OBc</i> | FNIFVGDLS | PEVTD | AT | LYQSF | -SVFS | SCS | DA | VMWDM | QKTGR | SRGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | |
| NAM8 | <i>S_pombe_NAM8</i> | YSIFVGDLS | PNVNE | FDVYSL | FASRYN | -SCKSA | KIMT | DPQTN | VS | RGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | | |
| | <i>A_fumigatus_NAM8</i> | YSIFVGDLS | PEVNE | YVLVSL | FQNRFP | -SCKSA | KIMT | DPIS | GMS | RGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | | |
| | <i>M_grisea_NAM8</i> | YSIFVGDLS | PEVNE | FVLVSL | FQARFP | -SCKSA | KIMT | DAMT | GQ | RGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | | |
| | <i>N_crassa_NAM8</i> | YSIFVGDLS | PEVNE | FVLVSL | FQSRFP | -SCKSA | KIMT | DAMT | GQ | RGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | | |
| | <i>Y_lipolytica_NAM8</i> | YSIFVGDLS | PEVNE | FVLVSL | FQSRFP | -SCKSA | KIMT | DAMT | GQ | RGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | | |
| | <i>D_hansenii_NAM8</i> | YSIFVGDLS | ASTTE | AHL | LAFFQKSF | PTSI | KTVR | VM | TPVSG | KRC | FGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | |
| | <i>A_gossypii_NAM8</i> | YSIFVGDLS | APNVT | ESQLF | ELFIS | RYST | LN | AKIV | FDQGT | CVS | KGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | |
| NGR1 | <i>K_lactis_NAM8</i> | FTLFVGDLS | APNVT | EAQLF | ELFIS | RYST | LN | AKIV | FDQGT | CVS | KGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | |
| | <i>C_glabrata_NAM8</i> | NSIFVGDLS | APSVT | EAQLF | ELFIS | RYST | LN | AKIV | FDQGT | CVS | KGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | |
| | <i>S_cerevisiae_NAM8</i> | CSIFVGDLS | APNVT | ESQLF | ELFIS | RYST | LN | AKIV | FDQGT | CVS | KGYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | |
| | <i>S_pombe_CSX1</i> | FSIFVGDLS | PTTED | SDLFMT | FRS | IYP | -SCTSA | KIVD | PVTGL | RKYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | | | |
| | <i>D_hansenii_NGR1</i> | YSIFVGDLS | GS | VT | EPMLF | ECFN | KVYP | NQV | KQAKIM | FD | PVTGL | RKYGF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ | | |
| | <i>A_gossypii_NGR1</i> | FSLFVGDLS | PTATE | AHL | LLSLF | QTKFK | -SVKT | VR | VM | TD | ITGAS | R | CF | GF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ |
| | <i>K_lactis_NGR1</i> | FSLFVGDLS | PTATE | AHL | LLSLF | QTKFK | -SVKT | VR | VM | TD | ITGAS | R | CF | GF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ |
| NGR1 | <i>C_glabrata_NGR1</i> | YSLFVGDLS | PTATE | AHL | LLSLF | QTKFK | -SVKT | VR | VM | TD | ITGAS | R | CF | GF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ |
| | <i>S_cerevisiae_NGR1</i> | FSLFVGDLS | PTATE | AHL | LLSLF | QTKFK | -SVKT | VR | VM | TD | ITGAS | R | CF | GF | FNK | WDAENA | IQQMNG | EWLGS | RA | IRVNW | ANQ |

a Recruiting function: RS-domain dependent



b Antagonist function: RS-domain independent



Exonic Splicing Silencers (ESSs)

Nature Reviews | **Genetics**

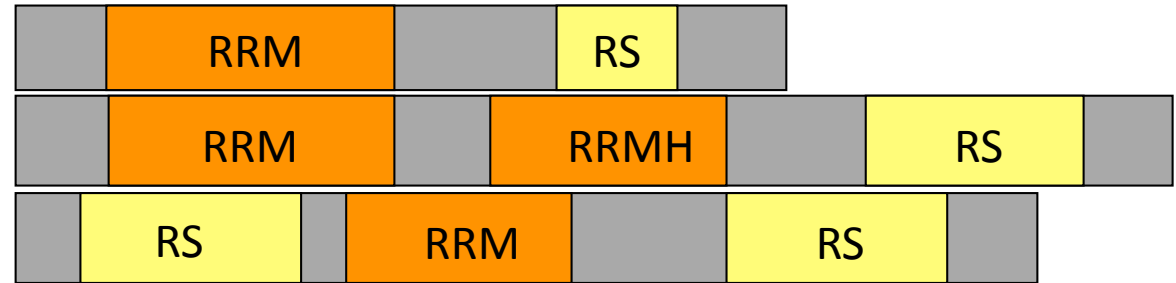
Domain organization of splicing factors

SR proteins

SC35, SR46, SRp20, 9G8, p54, SR86

SRp40, SRp55, SRp75, ASF, SR30C

TRA2, RNPS1



hnRNPs

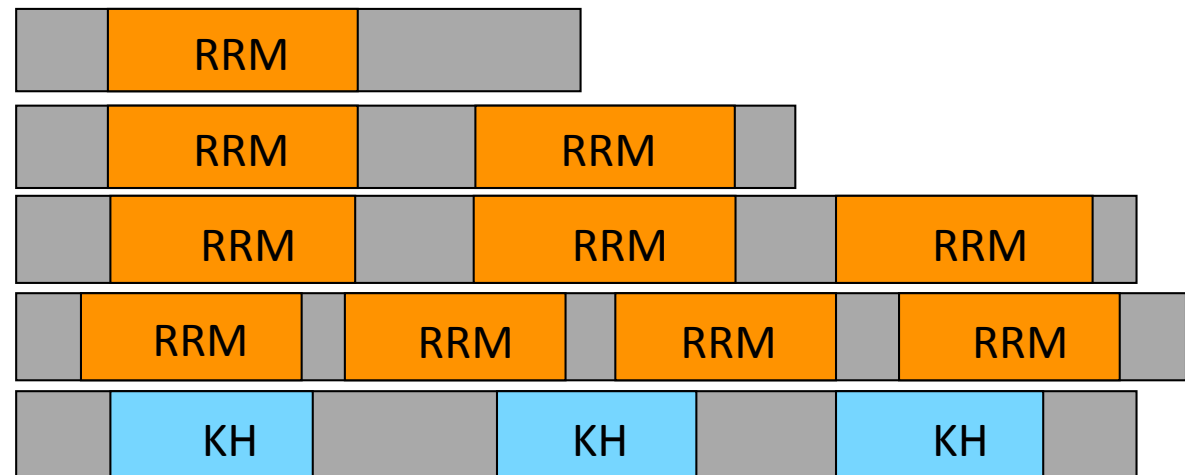
hnRNP-C, hnRNP-D, hnRNP-G

hnRNP-A, hnRNP-D, hnRNP-F/H, musashi

hnRNP-L, hnRNP-M, hnRNP-R

hnRNP-I

hnRNP-E



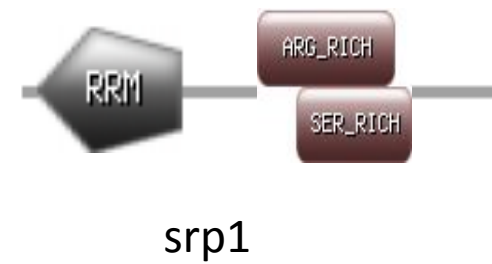
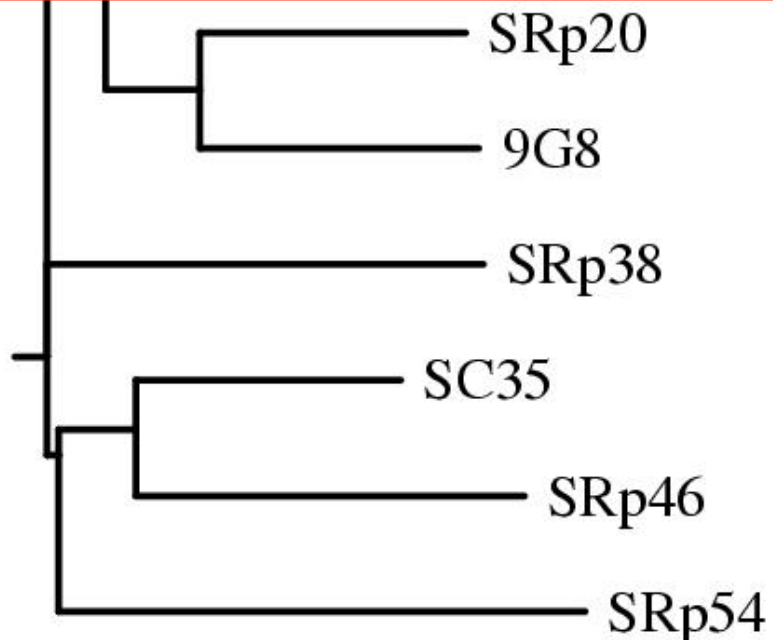
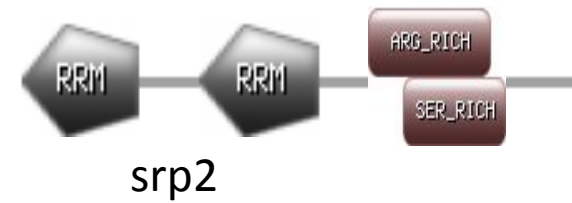
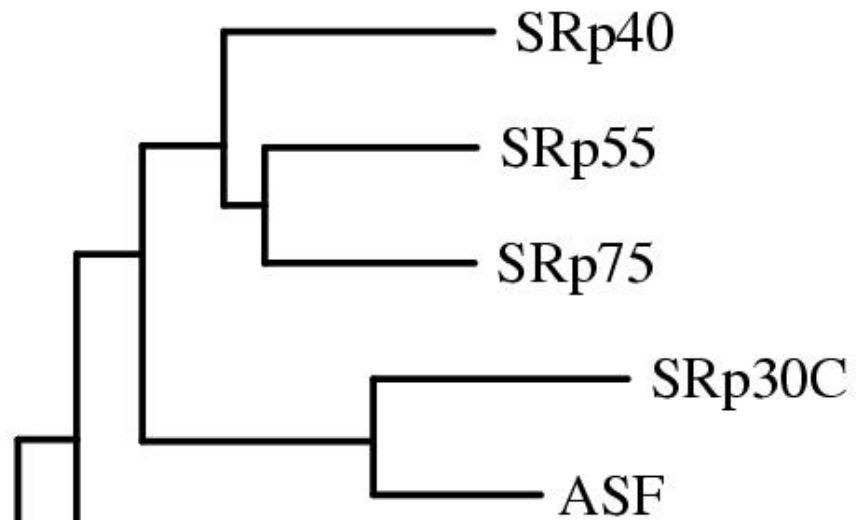
RRM: RNA Recognition Motif

RRMH: RRM homologous (similar to RRM)

RS: variable Arginine/Serine rich region (contains RS or SR dipetides)

KH: K-Homology domain

Human SR proteins



Starting hypotheses

Proteins with the same evolutionary origin (orthologs) are conserved
(have conserved functional sites)

Proteins with similar functions and related evolution (homologs) are also conserved

Proteins that perform similar interactions (Protein-RNA) but have opposite function
Are they conserved??

Can we model these differences/similarities with a simple model?

Objective

Build a Naïve-Bayes model with information from three types of data:

RRMs from enhancer proteins: SR proteins

RRMs from repressor proteins: hnRNP

KH-domains from repressor proteins: hnRNP

RRMs from other splicing factors (non-SR, non-hnRNP)

Can this model correctly classify a list of unknown domains?