

Title	Position-dependent termination and widespread obligatory frameshifting in Euplotes translation
Authors	Lobanov, Alexei V.;Heaphy, Stephen M.;Turanov, Anton A.;Gerashchenko, Maxim V.;Pucciarelli, Sandra;Devaraj, Raghul R.;Xie, Fang;Petyuk, Vladislav A.;Smith, Richard D.;Klobutcher, Lawrence A.;Atkins, John F.;Miceli, Cristina;Hatfield, Dolph L.;Baranov, Pavel V.;Gladyshev, Vadim N.
Publication date	2017
Original Citation	Lobanov, A. V., Heaphy, S. M., Turanov, A. A., Gerashchenko, M. V., Pucciarelli, S., Devaraj, R. R., Xie, F., Petyuk, V. A., Smith, R. D., Klobutcher, L. A., Atkins, J. F., Miceli, C., Hatfield, D. L., Baranov, P. V. and Gladyshev, V. N. (2016) 'Position-dependent termination and widespread obligatory frameshifting in Euplotes translation', Nature Structural and Molecular Biology, 24, pp. 61–68. doi: 10.1038/nsmb.3330
Type of publication	Article (peer-reviewed)
Link to publisher's version	https://www.nature.com/articles/nsmb.3330 - 10.1038/nsmb.3330
Rights	© 2017, Nature America, Inc., part of Springer Nature. All rights reserved. This is the peer reviewed version of the following article: Lobanov, A. V., Heaphy, S. M., Turanov, A. A., Gerashchenko, M. V., Pucciarelli, S., Devaraj, R. R., Xie, F., Petyuk, V. A., Smith, R. D., Klobutcher, L. A., Atkins, J. F., Miceli, C., Hatfield, D. L., Baranov, P. V. and Gladyshev, V. N. (2016) 'Position- dependent termination and widespread obligatory frameshifting in Euplotes translation', Nature Structural &Amp Molecular Biology, 24, pp. 61–68. doi: 10.1038/nsmb.3330, which has been published in final form at https://doi.org/10.1038/nsmb.3330.
Download date	2024-05-15 04:56:00
Item downloaded from	https://hdl.handle.net/10468/6518



University College Cork, Ireland Coláiste na hOllscoile Corcaigh



000-0

5



Supplementary Figure 1

Features of Euplotes genomes.

(a) Comparison Euplotes genomes in comparison with the genomes of other representative eukaryotes. The tree was constructed based on the sequences of 18S rRNA genes, and archaeal 16S rRNA gene (from Pyrococcus furiosis) was used as an outgroup. *number of contigs with telomeric repeats at both ends. (b) Distribution of telomeric repeat lengths in E. crassus (red) and E. focardii (black) macronuclear genomes. The X axis indicates the observed telomeric repeat number and the Y axis their frequencies. As expected, Euplotes genomes consist of gene-sized chromosomes capped by telomeres. The length of terminal repeats slightly varies; however, most chromosomes in both organisms have a double-stranded telomere length of 3.5 repeats (c) Sequence logo of subtelomeric regions at the 3' end of E. crassus nanochromosomes. 1000 randomly selected chromosome sequences with telomeric repeat GGGGTTTTGGGGTTTTGGGGTTTTGGGGG were chosen for constructing the logo. The logo detects a conserved positionspecific sequence motif associated with telomeric repeats. Abundance of high-quality telomeric sequences allowed an unbiased screen for motifs and patterns associated with telomere function. A previously described TCAA motif (Baird S. E. & Klobutcher L. A., Genes Dev 3, 585-597, 1989; Klobutcher, L. A. et al., Proc Natl Acad Sci USA 78, 3015-3019 ,1981) was readily detected with Weblogo (Crooks, G. E. et al, Genome Res 14, 1188-1190, 2004) in the subtelomeric region due to its conserved position relative to the telomere repeats. An analysis of sequences in the vicinity of telomeres with a pattern discovery suite MEME (Bailey, T. L. et al., Nucl Acids Res 34, W369-373, 2006) did not reveal additional common motifs.



Supplementary Figure 2

Features of the Euplotes transcriptome.

(a) Euplotes Sec and Cys tRNAs that decode TGA codons. Cys tRNA with the GCA anticodon and mitochondrial Trp tRNA with TCA anticodon are shown for comparison. In total we identified 183 tRNA genes in *E. crassus* and 337 genes in *E. focardii* based on their genomes analysis. (b) Frequency of introns of different lengths. The X axis indicates the length of introns in nucleotides, and the Y axis shows how many times they are found in the transcriptomes (log scale). Short introns (~25 nucleotides) is a characteristic feature of *Euplotes* transcriptomes. (c) Frequency of chromosomes with different numbers of RNA molecules transcribed from them. The X axis shows a number of transcripts per chromosome, and the Y axis how many such chromosomes are found in the genome. (d) E. crassus splice sites. Nucleotide conservation around exon-intron junction and intron-exon junctions. E. crassus. Transcriptomes were assembled *de novo* using Trinity (Haas, B. J. *et al., Nature Protoc*, 8, 1494-1512, 2013); no genomic template was used for the assembly of the transcriptome to ensure independence of the analysis. The assembly procedure produced 33,701 unique transcripts with an average length of 573 nucleotides in *E. crassus*. We obtained the *E. focardii* RNA-seq reads from (Keeling, P. J. *et al., PLoS Biol*, 12, e1001889, 2014).; this assembly produced 28,869 unique transcripts with an average length of 667 nucleotides. To identify introns we carried out pairwise alignments between the genome and the transcriptome for each species using FASTA (Pearson, W. *Curr Protoc Bioinf*, Chapter 3, Unit3 9, 2004) In total, we identified 21,798 introns in *E. crassus* and 18,747 in *E. focardii*. The most

frequent intron length was 25 nucleotides in both *E. crassus* and *E. focardii* with 2,895 and 2,631 occurrences, respectively. Using 10,000 intron sequences from E. crassus, we characterized sequence features of the exon-intron donor and intron-exon acceptor sites. We further aligned 32,350 E. crassus transcripts or their fragments (96%) to 18,032 genomic contigs, and similarly aligned 21,233 E. focardii transcripts (74%) to 16,950 genomic contigs. The majority of chromosomes had a single transcript aligning to them, 10,495 in *E. crassus* and 14,082 in *E. focardii*. Some chromosomes contained two or more predicted transcripts, which could be, at least in part, due to insufficient sequence coverage. Low coverage can result in missassembly of a single transcript as two or more, when reads matching internal positions are missing.



Nature Structural & Molecular Biology: doi:10.1038/nsmb.3330

Supplementary Figure 3

Termination at AAATAA and two mRNAs with long 3' UTRs.

In each panel ribosome footprints (top) and mRNA-seq reads (middle) are shown for a transcript whose ORF organization is shown at the bottom (red lines correspond to stop codons, and green lines to ATG codons). Identity of stop codons and adjacent 5' codons is indicated for the site of termination. Translated segments of ORFs are highlighted in blue. (a) An example of mRNA with termination at AAATAA. (b) mRNA of selenoprotein P22. The position of UGA Sec codon is shown in dark blue. (c) A single detected example of an mRNA with a long 3'UTR not containing SECIS structure.



Supplementary Figure 4

Metagene analysis of RNA-seq density surrounding frameshifting sites.

First nucleotide of a stop codon is shown as a zero coordinate. Only minor alteration of density associated with sequencing biases at specific nucleotides of frameshift sites can be seen.

Species	Assembler	Assembly size, kbp	Number of contias	Number of nanochromosomes*
Euplotes focardii	ABYSS NEWBLER	91,569 94 015	363,689 109 492	7,199 12 922
	SOAP	200,640	1,144,956	4
	SSAKE	118,465	374,877	8,879
	VELVET	114,730	301,971	4,996
Euplotes crassus	CELERA	19,350	12,326	247
	NEWBLER	59,563	56,588	14,194
	PCAP	64,474	70,328	8,097

Supplementary Table S1. E. crassus and E. focardii genome assemblies.

* Contigs containing both telomeric caps were designated as nanochromosomes. The assemblies shown in bold were used for further analyses.

Supplementary Note 1. *E. crassus* proteins with recoded and frameshift sites identified by mass spectrometry analyses.

a. Five out of nine selenoproteins (encoded by genes with UGA codon reassigned to code for selenocysteine) were detected by whole lysate high-throughput MS/MS analysis. Selenocysteine is shown in red. Sequences of the identified peptides are highlighted in yellow.

>eTR1

MDYSDTPQEESTHSYDYDLFVIGGGSGGLACAKVAQEAGAKVAVADFVKPTPKGTKW<mark>KVGGTCVNVGCIPK</mark>KLMHYSALLGNSYHDQVE SGWEHEKPSHDWG<mark>KMITNVNNHIRG</mark>INFGYKADMRKRGIKFHEKFASFVDPHTVQLVDKKGKTEMITSNYFVIATGGRPLYPDIPGAKE HAITSDDIFWMKDNPGKTLVVGASYVALEcAGFLHHFGNEVSVCVRSIFLRGFDQDMAQKIA<mark>KDMELSGINFIRDSVPTKIEKDEET</mark>GK LTcFLTVGGEETTVEVDTVLFAIGRYAVTADLNLGNAGLIAEKNGKFITDKYQ<mark>KTNVDNIYAIGDVLHGKLELTPTAIQAGRL</mark>LADRLF AGGTTTMDFYDVPTTIFTPLEYGcVGYSEEDAREEYGDFI<mark>KVYHTYFQPLEWNFAKS</mark>IYKERNcYVKIIVNTADNDRVIGFHILCPNAG EITQGIAIAIKVGVTKPQLDNCVGIHPTIAEEMTNLHIDKADNPDPIKSDC**U**S

>eP22

MESSDDKVGCVQSILVVLEGLNDDSSIITGLEILIKLIKNILKSPHEEKFRNIKKTNKAISTKLLSLSGIEDLILALGYKDDNDEFYVF DIDKYSDLYKLKRAIQEFHDEKRKKYMTPEELEKFEILQEQKRKFYEDNKKKAKARKDLENGMKFDREEKNQEEIKSS<mark>KANHLNFGANV</mark> VKFQPPAPASR**U**G

>eSelW2

MDSTTKGHIVVNYCGG**U**GYLPKA<mark>RYVQEAVENRFPGDFSFDLKA</mark>DVGKTG<mark>RLEVTVFVGDDTEGKLVHSKDKGQGFVKDSNVDSVLDSI</mark> AALLE

>eGPx1

MGAALCFKKRKE<mark>KLETTVESLFEISAEDIDGQEHLLADLAKD</mark>KKCIMVVNVASK**U**GLTKTHYTQMVKIHNKYKDKGFEIFAFPCNQFLS QEPGSNEDIKKFAREKYGAEFQLFS<mark>KIDVNGPNTHEVFRF</mark>CRRHSPLYDDETDTIQNIPWNFAKFLIDEEGNVVNYYSPKSNPDVCVPM IEEMLGL

>eGPx2

MGQVFFKSKKEKLATTVKSLFEISAKDIDGQTHLLADLAEGRKCTMVVNVASK**U**GLTKTHYKQMVKIHNK<mark>YRDHGFEIFAFPCNQFMSQ</mark> EPGTHEQIKK</mark>FAQE<mark>KYGAEFPLFSKVDVNGPDTHEVFKF</mark>CR<mark>RHSPLYDAEKDVVQNIPWNFAKF</mark>LIDE<mark>KGQVVEYYTPKQNPDLCVPKI EEM</mark>LGL

b. Sequences of proteins predicted to contain frameshifting. Sites of frameshifting are shown with an exclamation point highlighted in red. Sequences of the identified peptides are highlighted in yellow.

>comp7880_c0_seq1 AAATAA

MDNIPDYLVLRLNGTSFLDRREEILSIEYNSIFTFAEFKMEACRTLRVKYDPKCRLFDKEGIETFEDDLNMLKSHDVLYLASRGEDFDY SSVLNDYDRKDVLGEGGFGKVYHAVNRETGEDVAIKFMDISHYLTHADQIEEIYREADALQKLNHSHIISLHKAFVQRKEVILIMEYAG GGELKDRVEEMKDMDEIYARFIFQQICSAMSYCHNRGLIHRDLKLENVLFKEKGGMIIKIVDFGIAGVCKPQEKEKTDSGTLSYMPPEV LSGEKLEAGPGIDIWALGVMLYTMIYGKLPFYGDTEDEIINCIIKKKPSFKDKKTISKELKDLLIKVLNKDSDKRLSMFDLQNHKWME<mark>M QDEEILKSIEESKLEQEQEEEKK NE</mark>EDELIAFDKLNIKDDKKSSKAGSDYNSLSAHSSNPSGGRKKKKMRGSSPRSGMNGTGKKKKKT IKKKAT

>comp8353_c0_seq1 AAATAA

KKAHAVFNGGVTALCFSKKRTTLYTAGGDGTFLVWPVGAKPNPNQSVEPADFFSSPDLSNIPQIDNETDDSVKYYKELLQEQFLDSEIP RKKEFKAEISKELEDIRSHIVDLVEDNRKAQEIEQLERHEFVIDVKKKQSMEQDGWQQRQLITKQAKRKQLEYECLKERVKSATWSTME THSTACISLDSDLLLYNYGIRLRTPLEKKTLNRVLHFRRMELRQQITGMETRANKILDQSLFSNHDETYIMNRIRGVQHYEVDEEAPII TQAKSATKRKNKIAQQEASKSTADGALKKGKRKPKVDLNQFRLGANKPKLLDEDDFDDKLKDRAQNRDESNIAELRWKIVTKKKELEDL KKDIHILGSWDLLYEPSDLYTDGRKKMQIEMLQDVVFALKQEYNKEFERFQRFKDDQIFAIQERSNRITEILEDLKREEELFHPRTHPL ETPASILEVKPEEVTVQKYLTAEERAVEDEKHRIEQERLKALEGDNVGQRGIKNMLGGTYELKKNKGIMEETLEREEWMSKPIEDMTED EKLKFKEFQQREKELQEEKDKKRKAWDQELKKNRIEIEEICFKFEDELKKIHKKRLYYDMRVYEQELYIIRLTLMLHENKEIKQEAALI AEKKDKLEQELVESKNTINNFQRMYEDFDSEFKASSAIAEQEKGLRDLFPNAPFRQILAFVRNGGKAANRARGFRGQTEENTLEQEALA

>comp7670 c0 seq1 AAATAG, AATTAA MDEETIEKYCQALDLSVSPDNDLRKQAEAFIIEGMETPGFIAAMLHISSNPDLNRDRKIDITQAAAIQFKNIVETHWKYKDDEYAKEMR EDGYKVIIIPDETKTYVKENILTAYINVHSEKVAKQFDFIVRCITKHDFPDKWPDLANKVKDYIESDDLYGSEMFVGLYTLKSICKRYE YEFDAKREPLNEIADILFPRLEAITTCVEGDNSDQGSRLKNLIGHCFYISNQISLCKRYLDPSMLDFIVKFNTSALEAEIDNSLTQPTE SIEEIDHRAESFQWKLKMTAMNFLFRIFQKFSNPQYVNETMKPIAEHCINNYAEGIINLANTLIAKAK<mark>S</mark>SIYIDRQVLSYCFKVVSTSI NQTSYREMIKPLIPEILTSHCVPAMLLTEKDTEDFEADPVEFIRKARDPNPNIYTARNSVLEMIRNVTQHKSNQDKGALPDFLESFFGF LLENLSECIKQDAPDFRIKDALLLCLGQIAPTLLMYDQFHDQLNQVLTGAVFQDLTSENELVKYRALWVYGQCSRVPMEDDHRLEVGKL LFQLMNDENTAVKITASTSLYKNLRNNSMKEAFKSELASILEAYLGLMDTIDNEELIAGLEEVVSLYEDCIGPYAIELCSK<mark>IVENFN</mark>K ITGKEQEEEEYGTMGMATSGLVVTIRSIINSCKGDPETLLKLEPVIFPVVVRSLSADGCEYLDEAMDCITAILNFTQSATERMWALFPH LIK<mark>IIVGGPEDEEGGYAFDYFTSMEDYFR</mark>SLIKYGHDGMLTKKIGNDPVMILLIKGIIKILQLVKEGDVNTNAYICIVIVETLLLEFPG KLDQLLPTFIKILCTELSNKEITKEFRLHALTLVAHCFIYNCTLTLGALTDLKVLVPVCQNFFSYLKKFSEVEHLRGLIYGITALLRMD

>comp5116 c0 seq1 AATTAA MLPKPTNNMEKIKQQEYYKYKIRKVFECIAK<mark>ESNHNNDITN</mark>KNEIAYILR</mark>YFSQFPSEAQVTDYVIQKIEDDEPNDFIKLTKFEPYLL DVIENNEFEPSPPEHLLAAFRVLDTDKGRIPIDVLNNLLTTEGIPFRKEEMDSFOEFALDKSOKFVYYEDYVAKLVEENDKHVEEFLKE

>comp7341 c0 seq1 AACTAA MNNMIDCNHSPSLSDESVGADGDKECFSGERLEGSNNLPEEGITSISPDSLVNKASTLLRALTLFTSFSEFDTQNSPPPKTSKLFNDFS LKLNNMKTCLSQAKTSDSEEKLPALKEEIESIKASIGEELALTPLGQALLWNLIEGDNKDSSMKIFDELDHRCQDIANLHEVIAQKDA EIQTLSKQIRKLAKFRRTSSLVSEETEDGDSASQSDSGSMTQLSRSSSLLNKLNLNTKSRLTLCLNKVRDLMLVKELKEPLQKINTLSF NPGLLALEDAKEFLKNCFPLEVASFHFNKDSLLRNDLEKFLDVLLRTNEYVTDEIVLSNFVIDQDSLVKILSNFKNKEVVSFNSCKMSL SNPPEFGDSLDGATLKHLYLNFCGDKSHGDWASNPAHFENLINGLSHSPDLKASLKDIWMEGSGLKKDKARDILDTFGFHSTKIWILYG

LISNYNGLYRKIYDKEKFIND

YPTFKPPINO

>comp7882 c0 seq1 AAATAA MNPKGSKREKRVKGSKNKSAKEAFLKKIDAAMKVYDYVDETKDVKGKSERLNAINELQNLLQDQKSVSQLIIPNLDSCMQMIEKNIFRC LPNIKKSNLAFSETGIDQEEETDPAWPHVQGVYEFFLQLIMNDSIEVKLLKGYVTPEFVSRFLELFDSEEAVERDYLKNILHKLYAKLV PRRKMIRKAINETFYQLIHEGHKFNGASELLDILASIISGFAVPLREEHVIFFNNIIIRLHKVQTCSEFFEQLLRCSMLFLTKDKSLAI SLLKGLLKYWPFANCVKETLFLTELQEVLEIVDDDKIGDLVIPLFRRIVKCIGGTHLQVADRAMCFFENDYFLTKLRIYKDVTFPMLVP VIVELSENHWHKILQESLVALKVILKEIDSAAFDEAQQISK KDHRRFIVKPNVEKRTELDAKWERLNTTLKSTSAGFTPPDVPFKTSE

>comp6054 c0 seq1 AAATAA MEDEKNESIQNFKAMAECDDDGIAFQYLDSNNWDLAQAYDQYQNTHQFNQSNPTSTPAPTSFPGGADVDMSPGADAEESAFDDIPDIPN IGIPQMDQPSPVPQESNAPGSGLGNITSQFSNFASSIQSNLQNLTGGMFSGVMGGGMMGTDMNTQSNFSNRNLTAAQEFLFQFRKKNGM HVILPKFVNNTFEEIGQESKRLRRPVFFYLHNDKGDSCNIVDQSVIGEEMTRMLLNKYICVGVNVNTEEGRKLLTALEIPKAPFIGITY IDENGTLQNIGSRSGDEINVMALSEMDEAASGVFNAIFDGDTTDLTFHIEDSNLQLLETEEFKAEISAQMGNRTFDGYNEEPPRRRGPE IDPTTGFPVGMTPQQIQDKILKDQQRQEYAEIDEKNKVLIEER<mark>KK</mark>K<u>KQEENNLKR</u>KEELEKKAKIEKLEEEKKEMAEIVRSNLPEEPSE GTPDTITIQFRFPDGNHKQVRRFYKTDKVQLLYDYITSFGNENGFEAAHTHFSIIQNFPKKFFEDMNKTLEEEGLSNCTLMIKEHSHVE

FGLIFSDKYPILKTFFKFLAEKEVTHLTLDQWDSTYDLIRENPENLDNYDEYAAWPTLMDDFYQWYGENK

KFHMLKIIVDGITKKIRVSIIKKDFT >comp2566 c0 seq1 AAATAA MWGRKKKAEPKKSETKKRAPAKKATGTRKTKPPAKAKFSKLESKEEVKEEIKHDSETIFKVYATDOEMGRPIIGSEGIONLASDLKLDI ASSAELIVFMWHCECEEYGQISKSEFQKGCDK<mark>LGVKDFSHFK</mark>KSVPKKLSATLAMQDTPKEFRPFYKFAFTFHRTDGKNVPVETCQVL

>comp7194 c0 seq1 AAATAG MDTDLIDQIQKNMDQDPQLKDLFEPSSEENSQNDHGFTQGSKKFYSQDSAMAPPKVSRSKER<mark>ELAFLK<mark>R</mark>AQEIGLEPYNEYHGKKK</mark>TM IKKQTPGKDLIFSNIRHKESTSQQRRDHASRDEQQLSNKSLALKTGTNKEKIQKHKLHQSQRTVGKKTFECKEAKPSNGVAQKRVEVIE ISSKTSSSGNYSTPIESCPPIENCPSVESCPPVDHKTHEVVSLDSSNSDDKNIDDQPLNPKQKKRDKKKEQVDAKNARDFDSNPPKKSR MVSSTPTVNSEVMDKYLQQTPPDQIEIVEFDSRPKWPTENCSNFQDQLKLMASQRKRNKADTLGSNMFQEMKSTLSNIEASMDSPQGPI QKEYIHLKESIYPFWQSTFHHLEWNDDSDANKIKSREELKERMLDSLQYFSGHKLYRYADPADVKRGLLQNYPFVEDSDSKTEVKLLEG

LLCKLDPYIIVDENAVKAKFEQENVKEEYSYDRDKIVNLTPGEFDTLVQERENRNKIDKERKGMEQEIANLSGHKEFCEINANDLEEAY EDIKASHTEIESRMEKLKYNFEAVVYMLQGQVEVAQAPVATDYKDAILVNTGVIEDENKK<mark>VVQEGNTNVK</mark>KLEEITKFKRKLNHETWK NDKLKLEIKDLLERAIDVQLYKVTKDTQEIIKGNHRTKDEDEKKRLEDQINNLQENAGARIEVINKKKKKLRKEINEKRKENNELETRA RDLQTNVDQRNLIIDLR<mark>SKGPSGGDDRLQDPIKR</mark>FKEVATVRKYKEIVDQQKEEIEFLEDELERFRSRTFPSFANMHARQDYAD

LAQRIWNGEMDPDQMMDPNQMMDPNQ

>comp5973 c0 seq1 TTATAG MSHLKNFQFSSVQITEIDTYIEHLYSENMDLKLKGCISILYLCFSAENMEEMIEHESLLPAVSRILRDDYKKSLDLSLYLLNVFYAYSH FTEFHPLLIENQIGDTCVKIIEYEIKRYKARVNEYTKTAQLVKQTQQTPSADTDLKELQNNFRKEEKRLSVTIKKQEKVLFVTFHILLN LAEDLKIERKMKKRRIVPLLVSMLERNNPDLLYIVLSFLKKLSVFGSNKDDMLELDIMKKLNRFIPCQNALLTQTALRLLFNLSFDNEI RERVNAIGMIPKLVELLKVAQYRSILLRILYHLSSDDKIKATFAYTSCIPLVYQLVIHFPDAIIGKELIALAINLTTNKTNAALISQDD QLEALIERAFKYNDVLLFRVVRNIAQFGPVTNIDIYEKYMDKIIELTKQCGDNTDLQIELIGTLVYINIEKWDTVLSQGDFLDFIHNNL VSDYSEDDLVLETIMLIGTMCRSEKCAEAIAGSYIIGMLHELLGAKQEDDEMVQQILYTYHRLLYYRVTREIMLEQTQIVNVILELLND KNPNIRKLVNSTLDLVQLHDEIWKQEIKTKKFEMHNE<mark>VYL GLMEEYE</mark>AQAEALDEEALYDYYAQDPEALAALENGEFGEDDQWLDQND

VDESYHDCWKGMKRLFDPNDPDAGYKKYLSEHKN

>comp3853 c0 seq1 GTATAA MSEENKEEVKGTTHTDEDQYHHGFGNHFESEAIEGALPKHRNNPQQCKFGLYAEQISGTPFTYPRAKMQRSWLYRIMPTVAHPPYKALK DYNNLWIANFARDDDEEVFTTPQQMR<mark>WTPIDLPSEEITFVQGIQTV</mark>TGAGDPSMKAGINMGVYTCNTSMKNEAFFSSDGDIMIVPQLG KLSIMTEFGHIEAESWEVVVIPRGIKFAVEVNEDCRGYYCELYDGHLQIPDLGPIGTNGSANPRDFAIPKAKYFDETNEFRVIQKYLGK FFEYTIPHNIFDIVAWHGNYYPYKYDCHHFNTMGSISYDHPDPSVFTVLTCQTPDHGQAALDFAIFPPRWLSMEDTFRPPYFHRNTMNE FMGNVAGQYDAKEEGFSPGAVSLHSCMSAHGPEAEVVEK<mark>ASTCELKPQKVGE</mark>GCLAFMFETCYTMKVTKSFMHDLEGATDSYSVNSSKA

LEYDEVVDILEGKKNIGLGKEDKFKREMMEKIDRYIKKFQKYVGWT

>comp6951 c0 seq1 GAGTAA MYSTKFRRVMTMAPLLLANPALALCEEPSTADRIRGNYE<mark>NKIRFFAAPE</mark>KIFETFSNIREEDGQVYMSYQDFFHSLTPYNFVASKDDD DDDDDEENKDKEKEEREGYFDKFTPEIMTIVDANQDKKIDFNEYIFFITLLQLPEGEVMRIIEKVNPEERKINKAQFAKYLTKLRKCTA LGLKQMSKSFMPDGRKISTDEDHISKTILLHLFNDKEYITIEDFCELKSKLKHALLHYEFYQFDVDEDETISAESFAKSLLSCLNYTQA SKYSRRIHSLKLEGRVSFKEYVAFHNLIEKADIIKMKISTYRFLSLGMFRDLCDDFAKLDPYCNQNKVSISDTQIATFFK<mark>VLDEDENGA</mark>

FESLGTMIKEGWNNSEQVLKEIKAKEIENIKATIALESQKKIKELETTHQQEQLHYKKEMKKALEKTLAGLKMSYEDEIRLLKKNVKDQ ${\tt DKKITLLKKMCVRKDTQIQMLEEKAQSNEKQDKLQKEHRDILFELARAFKEGTTPGKDSTTNESATPY}$

>comp2483 c0 seq1 ATATAA, AAATAA MEKFGDLRASRHRVKHSKMKSKNHRQYEQEEVKHARSRPDRFDPPQIDEESKYSAGIDEAIQLVEITEKGVCKINPIAMNIVKGIKTKV GIISVVGPYRTGKSFLLNRLLGQQDGFEIGPTVQSCTRGIWIWGKPVKVSEDMHVILMDTEGLGSCNRTMNIDIKIFTLSVLLSSMFVY NCLNAIDENALEVLSLVVNLAKYISNQKKNDSMDVYNQANYSPYFMWVVRDFSLQMMPSEELEKAGHDPATYWDKLENQEAAAKEYLEK SLEAIDLGTINEENKRTVTKKNEIRKAIKNFFHQREATCLFRPINEEEKLRIVNKIPYEDLRKPFRKQVEHLINKIYYNVKPKSINGQT LTGKMFAQMLEEYTSSMNNNGMPEINTAWDRVMDTEIKRVLQESTTKINYQLQEVVIDKMPMPLKQLISIERNVRKSALKLLYDPNIKN APKDKLSRLQDKFIENLDEIFEGIFNENEIISKRQAKDLLPRMYQKIKAMINKGEFETIHDFSDIYGKMAISYFDNTNEPENYK<mark>I IQN</mark> FQINTVFEDLDEIMQTQVQRHESQNQEYETKLETKDHQIEHLNEQLKKEKTKNKDREQELRSKNMNIRSNLEEEIQMIK NQISNKDQQ

VRNQIPESALQELDQIDETLKETED

>comp6034 c0 seq2 ATATAA MSKNTKSKKQVTSNAKKGGNKKGKKAEPVQPPKEKKELAEWDLEDMPNFGFEPKKIAPTASKGPAVSGDKKKKGKKEKKTVEDTLISIE EAKRANPEEIARQETLITELKSQLEQKDQAIADLEKDQKEQFKQLTEQAQQLTEERDETRAALAVAEGQCNQKLDDFKQTVDRVNRNFL ENEKLI SELTSEKSNLKDIVFDLMFEKQKEGDSAPEGEEEITDEITDEIHGEFDRNTRRQAPQDNSQVKTLLDFANEQIKRLQTELKE

LPKPELGEIDEATEEKED

>comp8412 c0 seq1 ATATAA MSKNTKSKQEIDSTSKKLSRKERKNLEYIQYAKERKEYQKWEKEEADNLGFEGEESAPPVQNTSTAAGEKKKKGKKEKKPTDKVTTSPE EQKKAYSETVAKQDALIAELQVKLDQREKRMKDLKKTQEQKLSKLKDQTMKLTKERDDAKASLSKVEDKCNGRLGDYQELIE<mark>SVNRENL</mark> DNEKLI NDLTNDKANLKDIVFDLMFEKQKEGNSATEHPVVIPDNLQEEFNRKSQTKSSQNNPQLKTLLEFANEQIKRLQKEIKEARNH

VLQREGHLGLPSLEDDNEDDYFEPIG

>comp5528 c0 seq1 ATATAG MSSQEILANSITNTVDKEK<mark>SAQEEQDDEVIIDDQNPLLEDDLQI</mark>DEPEQKVNTDEPDQRNQEDEASENEQNLSDFINNTEFSYQSSST TQNLKNLLIQSTIGLALKLKPKLGKMLITNSDGGCCDDIRKSLDLNKQLLGEDVADLISVKQITWDESLQVSGTRYDYICITGSHFSQE FVQILEKIVPTVLSVVDDQERVFLLGPTSEEDISQFEDNVGDTIFESKKEEIDVPTKNSNSLGQFGDPDNHYSSENKDLIDEGIFDDDQ

EMPDVIKGSIQKIIESLIDLMRKYTRERILELRSKFEDKRNRWDEGTDEYNNLDAPFQKLSEWMEEYKDDAYSEDDDDDDDNFEEDDYLW SRSDSCYYKSCLEDKEAPLFFKETLEDFRENKEEVYRGIIELIPEDSQKLLEMIMERCEYMQSLQS

>comp4582 c0 seq1 TTATAA

LTVDSFILLADKKNCITLFSTFQDLISKIARKKHIFALNKNNEFAPNPMIGFIQNLCDKIYTITTDKEGKTAKEFLQDFEYCHNEEAEI DIPKPKIMKKMPPGIKKRLLADYNAKVEAAKKEVSKRAKNKVVISSSKTLLKEAFDLQDYHSFEYLFQYVEDKGIGYAELLHCLRNESN RKIFVLILDYVLTTLPEEEFE<mark>LIDVTNTTTQE</mark>ISLRELFPDIDLKEFCLALYDSKNVPGEIPLKSKYLYTKLEVYTKKYSTLDEKNRTK FLTTSLLVTGNTNPNEGYECLITKILDIISLYNIEIPIYYEGEVQEALRSNFKKLIFIRNYELVLKLQEIVKKNLRGLYEKL SQEHLS YISRLVQNSDAELGLEKDLLDSNEDVGVINSRQAVKDTLVFCLEQITLFDSYNQINFNDAPEKIIHMVKGFIHLGGFVNINVGFDSLDN KEVNEGDVSEIKRLVTVTEQYKEAKHKFSQRLLLTEFFSTFQKYNDLSYELIDVDSPLRWIIDCPGQESPNFFEYCIEQGNIDLAMKLI ESCDISEVISMFPLQERTISNLLNSPHIFEFLKKMSKEEEKIKKLIDRTNILEIPIMKLENSLSIDFDDEEDGNSGKPKYTQDQLTLYY FCYLKDSLVPKVGKEIPYFNLLFFNQGQFKYSLDELVKVLPLDKIKELNSLGYQIGKIISQKPVNTKDLIEI

>comp7073 c0 seq1 AAATAA, AAATAA, TCCTAA

KVSKESESLSQNKNKPIKRRKITEDDKVEHLLSNSNSNSQQQVNQVKPREEIKQPPQDPHKDQNMDADMARLESLDIPKVGRQPDKHKD HPMETDHDQKPQADANQQARDPEKPVRVPEDMRIPPSQPHVNPHLTEAPLRDAPSSQPLRAPQASPIHEIETAKKGKHVAPEVIRPDND VDMSKNMFENKSDRPKMQQERAQVVTTPQFTEQVPQRKDKAVVHKSISEIKKENDAPNHRDRKGRANDLSAK KLKFIDKYSTSQGRKE LGNMIRRISGPQVQGIVRLMRQFHVGNKEGKEFKFSLNTLTPAQCARVGMLIEGISDPGSASSTGRAQTGKPGSHGERSSSAVGSQDAS GAGRVSEREREIERRKRAEEEARYKERK KEHEMKLQERKKDELRRKEQEQRKEENRKFKEQQELLRRQEHERQQESDPHGPSYESKS PVPPTTSEQQEAARVKAHQEQLEQKRLEEEKRKQAEAERERIEQERRRAEEDKRRKAEELRKSEEQERQRELAKLRLEEERRKKKEAEE QRRREEERKRLELIKQKEEERRRQENSSPSKKSS KACEEERKR IREEEQRRRQQEEEDRKRQELQRKLKEDEERRLKAEQERKQREEEQRRIREEQERQRELAKKKEEERR KEEEERKRKEQEELRLREEQERKRREEERRRIEEERRR

Supplementary Note 2. Executable Analysis Document Supporting Proteomics Component.

1 Introduction

The vignette describes and reproduces all the steps that aimed to confirm frameshifts in the *Euplotes crassus* proteome. The global 8M urea soluble proteome was digested using conventional trypsin protocol and alternatively with Glu-C protease under high pH (7.5) conditions. The latter restricts specificity of Glu-C cleavages to C-terminal of glutamic acid (E). The peptides resulting from trypsin digest were fractionated using two different approaches: with strong cation exchange (SCX) and high pH reverse phase (HPRP) chromatographies. The peptides from Glu-C digest were fractionated using HPRP only.

The datasets were deposited to PRIDE and available by this link http://dx.doi.org/10.6019/PXD004333. Summary of the datasets shown in the table below:

Dataset Prefix	Digestion Enzyme	Fractionation Chromatography Type
Euplotes_1_SCX	trypsin	SCX
Euplotes_1_HPRP_1	trypsin	HPRP
Euplotes_1_HPRP_2	Glu-C (pH 7.5)	HPRP

Preprocessing of the raw files prior MS/MS searches was done in two steps. First, the raw files were processed with DeconMSn to correct for wrong assignments of monoisotopic peaks. The parameters are as follows:

```
DeconMSN.exe -I35 -G1 -F1 -L6810 -B200 -T5000 -M3 -XCDTA
```

At the second step the peak files were processed with DtaRefinery to perform post-acquisition recalibaration of parent ion mass-to-charge ratios. The peak lists (concatenated dta files in this case) were searched using MS-GF+ tool against 6-frame translated *Euplotes Crassus* genome concatenated with tentatively frameshifted sequences and common contaminants. The 6-frame translated FASTA file, DtaRefinery and MS-GF+ parameter files are available in extdata folder of the EuplotesCrassus.proteome package.

For example:

```
cat(readLines(fpath, n=12), sep = '\n')
## #Parent mass tolerance
## # Examples: 2.5Da or 30ppm
## # Use comma to set asymmetric values, for example "0.5Da,2.5Da" will set 0.5Da to the left (expMass<1
## PMTolerance=10ppm
##
## #Max Number of Modifications per peptide
## # If this value is large, the search will be slow
## NumMods=3
##
## #Modifications (see below for examples)
## StaticMod=C2H3N101, C, fix, any, Carbamidomethyl # Fixed Carbamidomethyl C (alkylation)
## StaticMod=C2H3N101, C, fix, any, Carbamidomethyl # Fixed Carbamidomethyl C (alkylation)
## Carbamidomethyl # Fixed Carbamidomethyl C (alkylation)
## StaticMod=C2H3N101, C, fix, any, Carbamidomethyl # Fixed Carbamidomethyl C (alkylation)
## StaticMod=C2H3N101, C, fix, any, Carbamidomethyl # Fixed Carbamidomethyl C (alkylation)
## StaticMod=C2H3N101, C, fix, any, Carbamidomethyl # Fixed Carbamidomethyl C (alkylation)
## StaticMod=C2H3N101, C, fix, any, Carbamidomethyl # Fixed Carbamidomethyl C (alkylation)
## StaticMod=C2H3N101, C, fix, any, Carbamidomethyl # Fixed Carbamidomethyl C (alkylation)
## StaticMod=C2H3N101, C, fix, any, Carbamidomethyl # Fixed Carbamidomethyl C (alkylation)
## StaticMod=C2H3N101, C, fix, any, Carbamidomethyl # Fixed Carbamidomethyl C (alkylation)
## StaticMod=C2H3N101, C, fix, any, Carbamidomethyl # StaticMod=C2H3N101, C, fix, any, Carbamidomethyl # Fixed Carbamidomethyl C (alkylation)
## StaticMod=C2H3N101, C, fix, any, Carbamidomethyl # Fixed Carbamidomethyl C (alkylation)
## StaticMod=C2H3N101, C, fix, any, Carbamidomethyl # StaticMod=C2H3N101, C, fix, any, Carb
```

2 Post MS/MS Search Analysis Steps

2.1 Prerequisites

2.1.1 Dowloading Datasets

To download the datasets we will take advantage of rpx R package. Note, this step may take awhile (10-30 min) depending on the speed of the internet connection. However, if they are downloaded the script will use the available datasets instead of downloading them again.

```
library(rpx)
id <- "PXD004333"
px <- PXDataset(id)
repoFiles <- pxfiles(px)
mzids <- grep('*msgfplus.mzid.gz', repoFiles, value=T)
system.time(pxget(px, mzids))
## user system elapsed
## 0.295 0.012 3.000</pre>
```

2.1.2 Reading Frameshift Marks

The FASTA files containing 595 sequences with frameshifts availabe as a part of this package and available as system.file("extdata", "Euplotes_Crassus_frameshifts.fasta", package="EuplotesCrassus.proteome"). There is an additional FASTA file with frameshift locations marked with exclamation mark !.

```
library(Biostrings)
fasta_clean <- readAAStringSet(
    system.file("extdata",
                     "Euplotes_Crassus_frameshifts.fasta",
                     package="EuplotesCrassus.proteome"),
    format="fasta", nrec=-1L, skip=0L, use.names=TRUE)
fasta_marks <- readAAStringSet(
    system.file("extdata",
                    "Euplotes_Crassus_frameshifts_with_mark.fasta",
                    package="EuplotesCrassus.proteome"),
    format="fasta", nrec=-1L, skip=0L, use.names=TRUE)
length(fasta_clean)</pre>
```

####

[1] 595

2.2 Processing of MS/MS Search Results

2.2.1 Trypsin Digest Fractionated by SCX

For processing of MS/MS identification we will use MSnID R package. First step is to read the LC-MS/MS datasets corresponding to 25 SCX fractions.

```
library(MSnID)
trypscx <- grep('Euplotes_1_SCX_.*msgfplus.mzid.gz', repoFiles, value=T)
trypscxPrj <- MSnID()
system.time(trypscxPrj <- read_mzIDs(trypscxPrj, trypscx, backend = 'mzR'))
## user system elapsed
## 4.829 0.214 5.106</pre>
```

Assess the peptide termini for their corresponding cleavage patterns. We will lleave peptides that resuted only from proper trypsin cleavave events. That is we won't allow peptide resulting from irregular clevages.

```
trypscxPrj <- assess_termini(trypscxPrj, validCleavagePattern="[KR]\\.[^P]")
trypscxPrj <- apply_filter(trypscxPrj, "numIrregCleavages == 0")</pre>
```

Note, that for this project we are interested only in peptides covering the sites of the frameshifting events. So if a peptide identification can be explained by a regular protein sequence we are not interested in pursuing this identification. The protein/accession names of normal (non-frameshifted) sequences starts with Contig or Contaminant. If the FASTA entry sequence is a results of the frameshift event if starts with comp. Therefore in the code below we retain only peptide-to-spectrum matches that can appear only due to frameshifted sequences.

Setting-up and optimizing filtering options for MS/MS identifications. Since the number of peptides mapping frameshifted sequences is rather low we will loosed up the FDR of the identification up to 5%, however, then follow-up with manual spectra validation.

```
trypscxPrj.fmsh$mme.ppm <- abs(mass_measurement_error(trypscxPrj.fmsh))
trypscxPrj.fmsh$score <- -log10(trypscxPrj.fmsh$`MS.GF.SpecEValue`)
trypscxPrj.fmsh <- apply_filter(trypscxPrj.fmsh, "mme.ppm < 10")</pre>
```

```
filtr <- MSnIDFilter(trypscxPrj.fmsh)
filtr$mme.ppm <- list(comparison="<", threshold=5.0)
filtr$score <- list(comparison=">", threshold=8.0)
```

```
#' pre-optimization with brute-force approach
filtr.grid <- optimize_filter(filtr, trypscxPrj.fmsh, fdr.max=0.05,</pre>
                              method="Grid", level="peptide", n.iter=20000)
evaluate_filter(trypscxPrj.fmsh, filtr.grid)
##
                        n
                   fdr
## PSM
            0.02970297 104
## peptide 0.03703704 56
## accession 0.04166667 50
#' fine tune with optimization using simulated annealing technique
filtr.sann <- optimize_filter(filtr.grid, trypscxPrj.fmsh, fdr.max=0.05,</pre>
                              method="SANN", level="peptide", n.iter=20000)
evaluate_filter(trypscxPrj.fmsh, filtr.sann)
##
                    fdr
                        n
           0.02941176 105
## PSM
## peptide 0.03636364 57
## accession 0.04081633 51
trypscxPrj.fmsh <- apply_filter(trypscxPrj.fmsh, filtr.sann)</pre>
show(trypscxPrj.fmsh)
## MSnID object
## Working directory: "."
## #Spectrum Files: 18
## #PSMs: 105 at 2.9 % FDR
## #peptides: 57 at 3.6 % FDR
## #accessions: 51 at 4.1 % FDR
```

Finally we will extract only those peptides that exactly span the frameshift sites. That is their sequences should be present/identifiable in normal FASTA file, however missing in the file with frameshifts masked with the exclamation mark !.

```
#' extract only those that map frameshift sites
library(dplyr)
pepSeq <- unique(trypscxPrj.fmsh$pepSeq)</pre>
pepSeqMapped_to_clean <- pepSeq %>%
    sapply(grep, x=fasta_clean) %>%
    sapply(length) %>%
    subset(.>0) %>%
    names
pepSeqMapped_to_with_marks <- pepSeq %>%
    sapply(grep, x=fasta_marks) %>%
    sapply(length) %>%
    subset(.>0) %>%
    names
pepSeqFmsh_trypscx <- setdiff(pepSeqMapped_to_clean, pepSeqMapped_to_with_marks)</pre>
print(pepSeqFmsh_trypscx)
## [1] "SAQEEQDDEVIIDDQNPLLEDDLQIDEPEQK" "WTPIDLPSEEITFVQGIQTVTGAGDPSMK"
## [3] "ESNHNNDITNKNEIAYILR"
                                          "KKKQEENNLKR"
```

Reporting extra information on the peptide sequences spanning frameshift sites: dataset, scan, charge, score, and mass measurement error.

```
meta_tryp_scx <- trypscxPrj.fmsh %>%
    apply_filter('pepSeq %in% pepSeqFmsh_trypscx') %>%
    psms %>%
```

####

spectrumFile	SpecEValue	MME (ppm)	spectrumID	charge	peptide
Euplotes_1_SCX_10_13Nov09_Falcon_09-09-14	3.41e-15	0.30	index=6106	3	K.SAQEEQDDEVIIDDQNPLLEDDLQIDEPEQK.V
Euplotes_1_SCX_10_13Nov09_Falcon_09-09-14	3.41e-15	0.30	index=6106	3	K.SAQEEQDDEVIIDDQNPLLEDDLQIDEPEQK.V
Euplotes_1_SCX_12_13Nov09_Falcon_09-09-14	1.53e-21	0.08	index=8908	2	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A
Euplotes_1_SCX_12_13Nov09_Falcon_09-09-14	1.07e-20	1.10	index=8896	2	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A
Euplotes_1_SCX_12_13Nov09_Falcon_09-09-14	7.29e-19	1.10	index=8897	2	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A
Euplotes_1_SCX_12_13Nov09_Falcon_09-09-14	2.17e-15	0.94	index=8895	3	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A
Euplotes_1_SCX_18_13Nov09_Falcon_09-09-15	9.27e-17	0.11	index=5912	2	K.ESNHNNDITNKNEIAYILR.Y
Euplotes_1_SCX_20_13Nov09_Falcon_09-09-15	2.23e-11	0.70	index=10317	3	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A
Euplotes_1_SCX_22_13Nov09_Falcon_09-09-15	4.36e-10	3.76	index=9720	3	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A
Euplotes_1_SCX_23_13Nov09_Falcon_09-09-15	2.47e-09	1.64	index=9440	3	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A
Euplotes_1_SCX_24_13Nov09_Falcon_09-09-15	3.42e-10	8.85	index=2127	3	R.KKKQEENNLKR.K

2.2.2 Trypsin Digest Fractionated by HPRP

All the processing steps are conceptually the same as in the section above.

```
tryphprp <- grep('Euplotes_1_HPRP_1_.*msgfplus.mzid.gz', repoFiles, value=T)</pre>
tryphprpPrj <- MSnID()</pre>
system.time(tryphprpPrj <- read_mzIDs(tryphprpPrj, tryphprp, backend = 'mzR'))</pre>
##
      user system elapsed
##
     2.716 0.175 2.945
tryphprpPrj <- assess_termini(tryphprpPrj, validCleavagePattern="[KR]\\.[^P]")</pre>
tryphprpPrj <- apply_filter(tryphprpPrj, "numIrregCleavages == 0")</pre>
tryphprpPrj.main <- apply_filter(tryphprpPrj, "!grepl('comp', accession)")</pre>
tryphprpPrj.fmsh <- apply_filter(tryphprpPrj, "grepl('comp', accession)")</pre>
tryphprpPrj.fmsh <- apply_filter(tryphprpPrj.fmsh,</pre>
                                 "!(peptide %in% peptides(tryphprpPrj.main))")
show(tryphprpPrj.fmsh)
## MSnID object
## Working directory: "."
## #Spectrum Files: 24
## #PSMs: 511 at 49 % FDR
## #peptides: 399 at 62 % FDR
## #accessions: 293 at 78 % FDR
tryphprpPrj.fmsh$mme.ppm <- abs(mass_measurement_error(tryphprpPrj.fmsh))</pre>
tryphprpPrj.fmsh$score <- -log10(tryphprpPrj.fmsh$`MS.GF.SpecEValue`)</pre>
tryphprpPrj.fmsh <- apply_filter(tryphprpPrj.fmsh, "mme.ppm < 10")</pre>
filtr <- MSnIDFilter(tryphprpPrj.fmsh)</pre>
filtr$mme.ppm <- list(comparison="<", threshold=5.0)</pre>
filtr$score <- list(comparison=">", threshold=8.0)
filtr.grid <- optimize_filter(filtr, tryphprpPrj.fmsh, fdr.max=0.05,</pre>
                               method="Grid", level="peptide", n.iter=20000)
evaluate_filter(tryphprpPrj.fmsh, filtr.grid)
##
                    fdr n
## PSM
           0.02631579 195
## peptide 0.04504505 116
## accession 0.07142857 75
filtr.sann <- optimize_filter(filtr.grid, tryphprpPrj.fmsh, fdr.max=0.05,</pre>
                               method="SANN", level="peptide", n.iter=20000)
evaluate_filter(tryphprpPrj.fmsh, filtr.sann)
##
                    fdr n
## PSM
             0.02604167 197
## peptide 0.04504505 116
## accession 0.07142857 75
tryphprpPrj.fmsh <- apply_filter(tryphprpPrj.fmsh, filtr.sann)</pre>
show(tryphprpPrj.fmsh)
## MSnID object
## Working directory: "."
## #Spectrum Files: 23
## #PSMs: 197 at 2.6 % FDR
## #peptides: 116 at 4.5 % FDR
```

```
#####
```

```
library(dplyr)
pepSeq <- unique(tryphprpPrj.fmsh$pepSeq)</pre>
pepSeqMapped_to_clean <- pepSeq %>%
    sapply(grep, x=fasta_clean) %>%
    sapply(length) %>%
    subset(.>0) %>%
   names
pepSeqMapped_to_with_marks <- pepSeq %>%
    sapply(grep, x=fasta_marks) %>%
    sapply(length) %>%
   subset(.>0) %>%
   names
pepSeqFmsh_tryphprp <- setdiff(pepSeqMapped_to_clean, pepSeqMapped_to_with_marks)</pre>
print(pepSeqFmsh_tryphprp)
## [1] "FFAAPEK"
                                         "ELAFLKRAQEIGLEPYNEYHGKKK"
## [3] "VVQEGNTNVKK"
                                         "WTPIDLPSEEITFVQGIQTVTGAGDPSMK"
## [5] "IIQNFQINTVFEDLDEIMQTQVQR"
                                        "KSSKACEEERRKR"
## [7] "LINDLTNDK"
                                        "LISELTSEK"
## [9] "IVENFNK"
                                        "LSQEHLSYISR"
## [11] "LINDLTNDKANLK"
meta_tryp_hprp <- tryphprpPrj.fmsh %>%
    apply_filter('pepSeq %in% pepSeqFmsh_tryphprp') %>%
   psms %>%
   select(spectrumFile,MS.GF.SpecEValue,mme.ppm,spectrumID,chargeState,peptide) %>%
   rename(SpecEValue = MS.GF.SpecEValue, charge = chargeState, `MME (ppm)`=mme.ppm) %>%
   mutate(spectrumFile = sub('_msgfplus.mzid.gz','',spectrumFile))
library(xtable)
print(xtable(meta_tryp_hprp, display = c('d','s','e','f','s','d','s')),
     include.rownames=FALSE,
      comment = FALSE,
      size='scriptsize',
      floating = F)
```

spectrumFile	SpecEValue	MME (ppm)	spectrumID	charge	peptide
Euplotes_1_HPRP_1_04_17Nov09_Falcon_09-09-14	7.58e-11	0.08	index=3031	1	R.FFAAPEK.I
Euplotes_1_HPRP_1_04_17Nov09_Falcon_09-09-14	2.44e-09	0.00	index=3046	2	R.FFAAPEK.I
Euplotes_1_HPRP_1_05_17Nov09_Falcon_09-09-14	1.46e-09	5.31	index=8245	3	R.ELAFLKRAQEIGLEPYNEYHGKKK.T
Euplotes_1_HPRP_1_06_17Nov09_Falcon_09-09-14	5.54e-10	2.21	index=759	2	K.VVQEGNTNVKK.L
Euplotes_1_HPRP_1_08_17Nov09_Falcon_09-09-14	5.93e-22	2.11	index=8644	2	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A
Euplotes_1_HPRP_1_08_17Nov09_Falcon_09-09-14	2.18e-21	0.78	index=8638	2	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A
Euplotes_1_HPRP_1_08_17Nov09_Falcon_09-09-14	3.05e-21	2.11	index=8646	2	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A
Euplotes_1_HPRP_1_08_17Nov09_Falcon_09-09-14	4.19e-16	0.82	index=8639	3	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A
Euplotes_1_HPRP_1_09_17Nov09_Falcon_09-09-14	1.19e-21	0.70	index=8806	2	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A
Euplotes_1_HPRP_1_09_17Nov09_Falcon_09-09-14	1.20e-21	1.57	index=8812	2	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A
Euplotes_1_HPRP_1_09_17Nov09_Falcon_09-09-14	5.49e-20	1.64	index=8802	2	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A
Euplotes_1_HPRP_1_09_17Nov09_Falcon_09-09-14	4.33e-15	1.53	index=8810	3	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A
Euplotes_1_HPRP_1_16_22Nov09_Falcon_09-09-14	4.51e-21	0.33	index=10684	2	K.IIQNFQINTVFEDLDEIMQTQVQR.H
Euplotes_1_HPRP_1_16_22Nov09_Falcon_09-09-14	1.36e-11	1.25	index=10678	3	K.IIQNFQINTVFEDLDEIMQTQVQR.H
Euplotes_1_HPRP_1_18_17Nov09_Falcon_09-09-15	5.08e-09	2.64	index=13785	2	K.KSSKACEEERRKR.E
Euplotes_1_HPRP_1_20_17Nov09_Falcon_09-09-15	1.91e-11	0.00	index=3425	1	K.LINDLTNDK.A
Euplotes_1_HPRP_1_22_17Nov09_Falcon_09-09-15	6.65e-11	1.67	index=3600	2	K.LISELTSEK.S
Euplotes_1_HPRP_1_22_17Nov09_Falcon_09-09-15	2.55e-10	0.78	index=3602	1	K.LISELTSEK.S
Euplotes_1_HPRP_1_22_17Nov09_Falcon_09-09-15	1.89e-09	0.49	index=2595	2	K.IVENFNK.I
Euplotes_1_HPRP_1_23_17Nov09_Falcon_09-09-15	3.01e-13	1.01	index=2200	2	K.LSQEHLSYISR.L
Euplotes_1_HPRP_1_24_17Nov09_Falcon_09-09-15	2.45e-16	1.41	index=2709	2	K.LINDLTNDKANLK.D

####

#accessions: 75 at 7.1 % FDR

2.2.3 Glu-C Digest Fractionated by HPRP

All the processing steps are conceptually the same as in the section above. The only substantial diffence is the specification of the enzyme digestion rule.

```
gluchprp <- grep('Euplotes_1_HPRP_2_.*msgfplus.mzid.gz', repoFiles, value=T)</pre>
gluchprpPrj <- MSnID()</pre>
system.time(gluchprpPrj <- read_mzIDs(gluchprpPrj, gluchprp, backend = 'mzR'))</pre>
##
     user system elapsed
##
     2.780 0.190 3.027
gluchprpPrj <- assess_termini(gluchprpPrj, validCleavagePattern="E\\.[^P]")</pre>
gluchprpPrj <- apply_filter(gluchprpPrj, "numIrregCleavages == 0")</pre>
gluchprpPrj.main <- apply_filter(gluchprpPrj, "!grepl('comp', accession)")</pre>
gluchprpPrj.fmsh <- apply_filter(gluchprpPrj, "grepl('comp', accession)")</pre>
gluchprpPrj.fmsh <- apply_filter(gluchprpPrj.fmsh,</pre>
                                 "!(peptide %in% peptides(gluchprpPrj.main))")
show(gluchprpPrj.fmsh)
## MSnID object
## Working directory: "."
## #Spectrum Files: 24
## #PSMs: 555 at 67 % FDR
## #peptides: 440 at 80 % FDR
## #accessions: 297 at 89 % FDR
gluchprpPrj.fmsh$mme.ppm <- abs(mass_measurement_error(gluchprpPrj.fmsh))</pre>
gluchprpPrj.fmsh$score <- -log10(gluchprpPrj.fmsh$`MS.GF.SpecEValue`)</pre>
gluchprpPrj.fmsh <- apply_filter(gluchprpPrj.fmsh, "mme.ppm < 10")</pre>
filtr <- MSnIDFilter(gluchprpPrj.fmsh)</pre>
filtr$mme.ppm <- list(comparison="<", threshold=5.0)</pre>
filtr$score <- list(comparison=">", threshold=8.0)
filtr.grid <- optimize_filter(filtr, gluchprpPrj.fmsh, fdr.max=0.05,</pre>
                               method="Grid", level="peptide", n.iter=20000)
evaluate_filter(gluchprpPrj.fmsh, filtr.grid)
##
                     fdr n
## PSM
            0.02222222 46
## peptide 0.03448276 30
## accession 0.05000000 21
filtr.sann <- optimize_filter(filtr.grid, gluchprpPrj.fmsh, fdr.max=0.05,</pre>
                               method="SANN", level="peptide", n.iter=20000)
evaluate_filter(gluchprpPrj.fmsh, filtr.sann)
##
                    fdr n
## PSM
             0.02222222 46
## peptide 0.03448276 30
## accession 0.05000000 21
gluchprpPrj.fmsh <- apply_filter(gluchprpPrj.fmsh, filtr.sann)</pre>
show(gluchprpPrj.fmsh)
## MSnID object
## Working directory: "."
## #Spectrum Files: 18
## #PSMs: 46 at 2.2 % FDR
```

## #peptides: 30 at 3.4 % FDR ## #accessions: 21 at 5 % FDR							
<pre>library(dplyr) pepSeq <- unique(gluchprpPrj.fmsh\$p pepSeqMapped_to_clean <- pepSeq %>% sapply(grep, x=fasta_clean) %>% sapply(length) %>% subset(.>0) %>% names pepSegMapped_to_with_marks <- pepSeg</pre>	epSeq)						
<pre>sapply(grep, x=fasta_marks) %>% sapply(length) %>% subset(.>0) %>% names</pre>	y 70~70						
<pre>pepSeqFmsh_gluchprp <- setdiff(pepS print(pepSeqFmsh_gluchprp)</pre>	eqMapped_	_to_clean,	pepSeqMapp	ped_to_	with_marks)		
<pre>## [1] "NFNKITGKEQEEEE" ## [3] "NLDNEKLINDLTNDKANLKDIVFDLMF ## [5] "MQDEEILKSIEESKLEQEQEEEKKNE"</pre>	E" '	'SVNRENLDNE 'NKIRFFAAPE 'VYLGLMEEYE	EKLINDLTNDA EKIFE" E"	ANLKDI	VFDLMFE"		
<pre>meta_gluc_hprp <- gluchprpPrj.fmsh %>% apply_filter('pepSeq %in% pepSeqFmsh_gluchprp') %>% psms %>% select(spectrumFile,MS.GF.SpecEValue,mme.ppm,spectrumID,chargeState,peptide) %>% rename(SpecEValue = MS.GF.SpecEValue, charge = chargeState, `MME (ppm)`=mme.ppm) %>% mutate(spectrumFile = sub('_msgfplus.mzid.gz','',spectrumFile))</pre>							
<pre>print(xtable(meta_gluc_hprp, display = c('d','s','e','f','s','d','s')),</pre>							
spectrumFile	SpecEValue	MME (ppm)	spectrumID	charge	peptide		
Euplotes_1_HPRP_2_06_22Nov09_Falcon_09-09-15 Euplotes_1_HPRP_2_08_25Nov09_Falcon_09-09-15 Euplotes_1_HPRP_2_09_25Nov09_Falcon_09-09-15 Euplotes_1_HPRP_2_09_17Nov09_Falcon_09-09-17 Euplotes_1_HPRP_2_09_17Nov09_Falcon_09-09-17 Euplotes_1_HPRP_2_12_17Nov09_Falcon_09-09-17 Euplotes_1_HPRP_2_12_17Nov09_Falcon_09-09-17 Euplotes_1_HPRP_2_12_17Nov09_Falcon_09-09-17 Euplotes_1_HPRP_2_215_17Nov09_Falcon_09-09-17 Euplotes_1_HPRP_2_222_22Nov09_Falcon_09-09-17 Euplotes_1_HPRP_2_16_17Nov09_Falcon_09-09-17	6.80e-07 3.78e-17 3.33e-07 5.74e-16 5.03e-07 2.09e-09 1.62e-07 2.83e-07 2.17e-07 2.12e-08	2.95 0.19 0.57 0.44 1.11 0.43 0.07 1.61 0.10 0.88	index=13369 index=9982 index=9974 index=10771 index=10770 index=3933 index=3930 index=1758 index=6671 index=66753	2 3 4 3 4 3 2 2 1 1	E.NFNKITGKEQEEEE.Y E.SVNRENLDNEKLINDLTNDKANLKDIVFDLMFF E.SVNRENLDNEKLINDLTNDKANLKDIVFDLMFF E.NLDNEKLINDLTNDKANLKDIVFDLMFE.K E.NLDNEKLINDLTNDKANLKDIVFDLMFE.K E.NKIRFFAAPEKIFE.T E.NKIRFFAAPEKIFE.T E.MQDEEILKSIEESKLEQEQEEEKKNE.E E.VYLGLMEEYE.A E.VYLGLMEEYE.A		

2.3 Compendium of Peptides Covering Frameshift Locations

Final set of peptides and corresponding references to LC-MS/MS datasets and spectra. Overall, **4**, **11**, and **6** unique peptide sequences spanning the frameshift sites were identified in trypsin/SCX, trypsin/HPRP, and 'Glu-C/HPRP' experiments, respectively.

spectrumFile	SpecEValue	MME (nnm)	spectrumID	charge	pentide	experiment
Euplotes 1 SCX 10 13Nov09 Ealcon 09-09-14	3.41e-15	0.30	index=6106	3	K.SAQEEQDDEVIIDDQNPLLEDDLQIDEPEQK.V	trypsin/SCX
Euplotes 1 SCX 10 13Nov09 Falcon 09-09-14	3.41e-15	0.30	index=6106	3	K.SAQEEQDDEVIIDDQNPLLEDDLQIDEPEQK.V	trypsin/SCX
Euplotes 1 SCX 12 13Nov09 Falcon 09-09-14	1.53e-21	0.08	index=8908	2	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A	trypsin/SCX
Euplotes 1 SCX 12 13Nov09 Falcon 09-09-14	1.07e-20	1.10	index=8896	2	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A	trypsin/SCX
Euplotes 1 SCX 12 13Nov09 Falcon 09-09-14	7.29e-19	1.10	index=8897	2	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A	trypsin/SCX
Euplotes 1 SCX 12 13Nov09 Falcon 09-09-14	2.17e-15	0.94	index=8895	3	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A	trypsin/SCX
Euplotes 1 SCX 18 13Nov09 Falcon 09-09-15	9.27e-17	0.11	index=5912	2	K.ESNHNNDITNKNEIAYILR.Y	trypsin/SCX
Euplotes_1_SCX_20_13Nov09_Falcon_09-09-15	2.23e-11	0.70	index=10317	3	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A	trypsin/SCX
Euplotes_1_SCX_22_13Nov09_Falcon_09-09-15	4.36e-10	3.76	index=9720	3	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A	trypsin/SCX
Euplotes_1_SCX_23_13Nov09_Falcon_09-09-15	2.47e-09	1.64	index=9440	3	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A	trypsin/SCX
Euplotes_1_SCX_24_13Nov09_Falcon_09-09-15	3.42e-10	8.85	index=2127	3	R.KKKQEENNLKR.K	trypsin/SCX
Euplotes_1_HPRP_1_04_17Nov09_Falcon_09-09-14	7.58e-11	0.08	index=3031	1	R.FFAAPEK.I	trypsin/HPRP
Euplotes_1_HPRP_1_04_17Nov09_Falcon_09-09-14	2.44e-09	0.00	index=3046	2	R.FFAAPEK.I	trypsin/HPRP
Euplotes_1_HPRP_1_05_17Nov09_Falcon_09-09-14	1.46e-09	5.31	index=8245	3	R.ELAFLKRAQEIGLEPYNEYHGKKK.T	trypsin/HPRP
Euplotes_1_HPRP_1_06_17Nov09_Falcon_09-09-14	5.54e-10	2.21	index=759	2	K.VVQEGNTNVKK.L	trypsin/HPRP
Euplotes_1_HPRP_1_08_17Nov09_Falcon_09-09-14	5.93e-22	2.11	index=8644	2	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A	trypsin/HPRP
Euplotes_1_HPRP_1_08_17Nov09_Falcon_09-09-14	2.18e-21	0.78	index=8638	2	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A	trypsin/HPRP
Euplotes_1_HPRP_1_08_17Nov09_Falcon_09-09-14	3.05e-21	2.11	index=8646	2	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A	trypsin/HPRP
Euplotes_1_HPRP_1_08_17Nov09_Falcon_09-09-14	4.19e-16	0.82	index=8639	3	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A	trypsin/HPRP
Euplotes_1_HPRP_1_09_17Nov09_Falcon_09-09-14	1.19e-21	0.70	index=8806	2	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A	trypsin/HPRP
Euplotes_1_HPRP_1_09_17Nov09_Falcon_09-09-14	1.20e-21	1.57	index=8812	2	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A	trypsin/HPRP
Euplotes_1_HPRP_1_09_17Nov09_Falcon_09-09-14	5.49e-20	1.64	index=8802	2	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A	trypsin/HPRP
Euplotes_1_HPRP_1_09_17Nov09_Falcon_09-09-14	4.33e-15	1.53	index=8810	3	R.WTPIDLPSEEITFVQGIQTVTGAGDPSMK.A	trypsin/HPRP
Euplotes_1_HPRP_1_16_22Nov09_Falcon_09-09-14	4.51e-21	0.33	index=10684	2	K.IIQNFQINTVFEDLDEIMQTQVQR.H	trypsin/HPRP
Euplotes_1_HPRP_1_16_22Nov09_Falcon_09-09-14	1.36e-11	1.25	index=10678	3	K.IIQNFQINTVFEDLDEIMQTQVQR.H	trypsin/HPRP
Euplotes_1_HPRP_1_18_17Nov09_Falcon_09-09-15	5.08e-09	2.64	index=13785	2	K.KSSKACEEERRKR.E	trypsin/HPRP
Euplotes_1_HPRP_1_20_17Nov09_Falcon_09-09-15	1.91e-11	0.00	index=3425	1	K.LINDLTNDK.A	trypsin/HPRP
Euplotes_1_HPRP_1_22_17Nov09_Falcon_09-09-15	6.65e-11	1.67	index=3600	2	K.LISELTSEK.S	trypsin/HPRP
Euplotes_1_HPRP_1_22_17Nov09_Falcon_09-09-15	2.55e-10	0.78	index=3602	1	K.LISELTSEK.S	trypsin/HPRP
Euplotes_1_HPRP_1_22_17Nov09_Falcon_09-09-15	1.89e-09	0.49	index=2595	2	K.IVENFNK.I	trypsin/HPRP
Euplotes_1_HPRP_1_23_17Nov09_Falcon_09-09-15	3.01e-13	1.01	index=2200	2	K.LSQEHLSYISR.L	trypsin/HPRP
Euplotes_1_HPRP_1_24_17Nov09_Falcon_09-09-15	2.45e-16	1.41	index=2709	2	K.LINDLTNDKANLK.D	trypsin/HPRP
Euplotes_1_HPRP_2_06_22Nov09_Falcon_09-09-15	6.80e-07	2.95	index=13369	2	E.NFNKITGKEQEEEE.Y	Glu-C/HPRP
Euplotes_1_HPRP_2_08_25Nov09_Falcon_09-09-15	3.78e-17	0.19	index=9982	3	E.SVNRENLDNEKLINDLTNDKANLKDIVFDLMFE.K	Glu-C/HPRP
Euplotes_1_HPRP_2_08_25Nov09_Falcon_09-09-15	3.33e-07	0.57	index=9974	4	E.SVNRENLDNEKLINDLTNDKANLKDIVFDLMFE.K	Glu-C/HPRP
Euplotes_1_HPRP_2_09_17Nov09_Falcon_09-09-17	5.74e-16	0.44	index=10771	3	E.NLDNEKLINDLTNDKANLKDIVFDLMFE.K	Glu-C/HPRP
Euplotes_1_HPRP_2_09_17Nov09_Falcon_09-09-17	5.03e-07	1.11	index=10770	4	E.NLDNEKLINDLTNDKANLKDIVFDLMFE.K	Glu-C/HPRP
Euplotes_1_HPRP_2_12_17Nov09_Falcon_09-09-17	2.09e-09	0.43	index=3933	3	E.NKIRFFAAPEKIFE.T	Glu-C/HPRP
Euplotes_1_HPRP_2_12_17Nov09_Falcon_09-09-17	1.62e-07	0.07	index=3930	2	E.NKIRFFAAPEKIFE.T	Glu-C/HPRP
Euplotes_1_HPRP_2_15_17Nov09_Falcon_09-09-17	2.83e-07	1.61	index=1758	2	E.MQDEEILKSIEESKLEQEQEEEKKNE.E	Glu-C/HPRP
Euplotes_1_HPRP_2_21_22Nov09_Falcon_09-09-17	2.17e-07	0.10	index=6671	1	E.VYLGLMEEYE.A	Glu-C/HPRP
Euplotes_1_HPRP_2_22_22Nov09_Falcon_09-09-17	2.12e-08	0.88	index=6753	1	E.VYLGLMEEYE.A	Glu-C/HPRP

##\$##

3 Manual Validation

Manual valiation was perfomed by LCMSSpectator. The spectra that have passed the consensus opinion of 5 independed experts are shown below. Necessary raw and mzldenML files to reproduce the analysis are available at http://dx.doi.org/10.6019/PXD004333. Note, the MS/MS scan number is not the same identifier as spectrumID in the table above.

SAQEEQDDEVIIDDQNPLLEDDLQIDEPEQK



SVNRENLDNEKLINDLTNDKANLKDIVFDLMFE



####



WTPIDLPSEEITFVQGIQTVTGAGDPSMK

ESNHNNDITNKNEIAYILR



20

FFAAPEK



#####

Nature Structural & Molecular Biology: doi:10.1038/nsmb.3330

IIQNFQINTVFEDLDEIMQTQVQR



21

LINDLTNDK



####

LINDLTNDKANLK



IVENFNK



LISELTSEK



LSQEHLSYISR



NKIRFFAAPEKIFE



VYLGLMEEYE



4 Session Information

All software and respective versions used in this document, as returned by sessionInfo() are detailed below.

- R version 3.2.4 (2016-03-10), x86_64-apple-darwin13.4.0
- Locale: en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
- Base packages: base, datasets, graphics, grDevices, methods, parallel, stats, stats4, utils
- Other packages: BiocGenerics 0.16.1, BiocStyle 1.8.0, Biostrings 2.38.4, dplyr 0.5.0, IRanges 2.4.8, knitr 1.12.3, MSnID 1.7.3, Rcpp 0.12.7, rpx 1.6.0, S4Vectors 0.8.11, xtable 1.8-2, XVector 0.10.0
- Loaded via a namespace (and not attached): affy 1.48.0, affyio 1.40.0, assertthat 0.1, Biobase 2.30.0, BiocInstaller 1.20.3, BiocParallel 1.4.3, bitops 1.0-6, chron 2.3-47, codetools 0.2-14, colorspace 1.2-6, data.table 1.9.6, DBI 0.5-1, digest 0.6.10, doParallel 1.0.10, evaluate 0.8.3, foreach 1.4.3, formatR 1.3, futile.logger 1.4.3, futile.options 1.0.0, ggplot2 2.1.0.9000, grid 3.2.4, gtable 0.2.0, highr 0.5.1, htmltools 0.3.5, impute 1.44.0, iterators 1.0.8, lambda.r 1.1.9, lattice 0.20-33, lazyeval 0.2.0, limma 3.26.9, magrittr 1.5, MALDIquant 1.14, MSnbase 1.18.1, munsell 0.4.3, mzID 1.8.0, mzR 2.4.1, pcaMethods 1.60.0, plyr 1.8.4, preprocessCore 1.32.0, ProtGenerics 1.2.1, R.cache 0.12.0, R.methodsS3 1.7.1, R.oo 1.20.0, R.utils 2.3.0, R6 2.1.2, RCurl 1.95-4.8, reshape2 1.4.1, rmarkdown 0.9.5, scales 0.4.0, stringi 1.1.1, stringr 1.1.0, tibble 1.2, tools 3.2.4, vsn 3.38.0, XML 3.98-1.4, yaml 2.1.13, zlibbioc 1.16.0

SUPPLEMENTARY NOTE 3. Representative profiles of ribosome density mapped



to *E. crasus* transcripts and supporting BLAST hits alignments.

Supplementary Note Figure 1. Supporting information for +1 frameshifting at AAT_TAA. Left panel: density of ribosome footprints (top) and mRNA-seq reads (middle) for a transcript whose ORF is shown at the bottom (red lines correspond to stop codons, and green lines to ATG codons). Identity of stop codons and adjacent 5' codons is indicated for the frameshift site and for the site of termination. Translated segments of ORFs are highlighted in blue. Right panel shows protein sequence produced with inferred frameshifting (top) and its alignment to the closest BLAST hit (bottom).



MLYGKT SVNTTL EFLRVI MLRIRN EVTVPS FHLNGL ILTDGD DIVQFV	EQIEN TEG <mark>RI</mark> QTEPI GDTQY SLHYL EDPEV IHDMP KFNDF	NLNPDFVTYFEMDYYFEXIQXIKVEVFDVDVTRLERIGNFETTIGEIMG ILRTEKVATSNDLIYFSLRINDLVSNKGWFGSDDPFIFIERARENDQE RNDLNPTMRVLKYEAKEICNGDLQCPLKFKVYSWRNSGHHKFFGEFETT NLFKDGAQQKSICSFEEFFIERASFFDFLHSGWKMNLMVCVDFTASNG NPTGEFNDYQNAIRQVGNILELYDYNRQYPCYGFGGIPRYSGSNQVSHC DGVNGILESVGFSLLNCGLYGPTNFGECMRKTVDVIKERMDERMYHTLL ITRDIIVEGSHYPLSIIIIGLGESSFDKMIELDGDDVVLKNTRGEATRR RHLSKQALAEEVLEEVPEQVVSYLSQNNIKLDEVN	
>emb C Length	DW786 =554	01.1 copine family protein [Stylonychia lemnae]	
Score Ident	= ities	420 bits (1080), Expect = 3e-137, Method: Compositional matri: = 211/516 (41%), Positives = 330/516 (64%), Gaps = 21/516 (4%	x adjus <mark>t</mark> .)
Query	16	EKITLHISCRKLADLDIITVSDPVCHVYIADSDHPDDWMLYGKTEQIENNLNPDFVTYFE	75
Sbjct	19	QRVSLSISCRNLKNLDVLSKSDPMCEVYIKDR-KTTNWTLLGKTETINNNLNPDFSSIIY	77
Query	76	MDYYFEKIQKIKVEVFDVDVTRLERIGNFETTLGEIMGSVNTTLTE	121
Sbjct	78	CDYFFEREQNIKFDLYDIDNQQHTSRDFIGSNETTLGGIIGSMQQTYVADLKDNKSTRSR	137
Query	122	GRIL-RTEKVATSNDLIVFSLRINDLVSNKGWFCGS-DDPFIFIERARE-NDQEEFLRVI	178
Sbjct	138	GKIVVRLDNVNTTNDEVRLRVSARVQSNAGCCGTQDNPYYIISRARDVNNHKDFVRVY	195
Query	179	QTEPIRNDLNPTWRYLKYEAKEICNGDLQCPLKFKVYSWRNSGHHKFFGEFETTMLRIRN	238
Sbjct	196	++ + N P W K + +ICNG P+KF++YS SG + +GE I++ ++++ KSSAMLNSTQPMWNVQKIKLSQICNGINNLPIKFELYSQNISGTDQAYGEGITSIEQLQS	255
Query	239	GDTQYNLFKDGAQQKSICSFEEFFIEERASFFDFLHSGWKMNLMVCVDFTASNGEVTVPS	298
Sbjct	256	G + + K + + F I E +F ++L SGW +N+ +D+TASNGE T P+ GQKSVEITDKKRKIKGSLNIDNFVIREMPNFMEYLRSGWAINMSFAIDYTASNGEKTDPN	315
Query	299	SLHYLNPTGE-FNDYQNAIRQVGNILELYDYNRQYPCYGFGGIPRYSGSNQVSHCFHLNG	357
Sbjct	316	SLH +P+G N Y+ A+ VG ++E Y N+ + +GFGGIPR++GSNQ+SHCF+LNG SLHKQDPSGRNLNQYEQALLSVGKVMEPYALNQMFATFGFGGIPRFTGSNQISHCFNLNG	375
Query	358	LEDPEVDGVNGILESYQFSLLNCGLYGPTNFGECMRKTVDYIKERMDERMYHILLILTDG	417
Sbjct	376	P++ G+ + +Y+ ++ GL GPT+F ++ + Y+++ + MYH L+I+TDG SVSPQIQGLQNVYMAYKRTIHQIGLAGPTHFSSVLQSLLVYVQQCLQFQMYHCLMIITDG	435
Query	418	DIHDMPITRDIIVEGSHYPLSIIIIGLGESSFDKMIELDGDDVVLKNTRGEATRRDIVOF	477
Sbict	436	+IHDMP T D+IVE S +P+SIIIIG+G F+KM LD D+ L+N++G+ RDIVQF EIHDMPATIDLIVELSRFPVSIIIIGVGNEGFEKMNFLDSDNOALRNSKGOVAARDIVOF	495

Supplementary Note Figure 2. Supporting information for +2 frameshifting at AGA_TAA. See Supplementary Fig. S8 for the legend.



>COMP99 MSRILEF LNLPSTF KVYQKEG LTTLPRT LSRFKSS	013_C0 RSGSVS RSTRST GFQNTH GRNFY SGTTVL	9_seq1_AGGTAA(+1) SNNPLNCSTAKQQFSFSKADRFGQTMTRGLSSGY TTFGYGDKINTNKRINDSPSPGTYKIPSGFEPGKKGKVYSFRCSWKSYS KASDPNVPGPGTYKSVPTFGNKGRKFTLKGKLKDLPSSVKVPGPGAYKD YSKFKSVCNGAIGPPTNTRFKDQTQTSAEIPGPGQYTPKTGLTGTGHYF LGGARRASTHRVSNDGPGPGSYILPSDFGYPTSFHRRRKKSRRATSQL	
>gb EJY Length:	(73919 =322	9.1 hypothetical protein OXYTRI_04827 [Oxytricha trifallax]	
Score Identi	= 1 Lties	170 bits (430), Expect = 1e-46, Method: Compositional matrix a = 111/289 (38%), Positives = 150/289 (52%), Gaps = 36/289 (129	adjust %)
Query	1	MSRILERSGSVSNNPLNCSTAKQQFSFSKADRFGQTMTRGLSSGYLNLPSTRSTRSTFG	60
Sbjct	1	MSQVVSTPQQINQHPVNNSTSKQLFSFSKAERFGQPKSMMNHRIAYDLPSMKSTRAAGLG	60
Query	61	YGDKIMTNKRNDSPSPGTYKIPSGFEPGKKGKVYSFRCSWKSYSKVYQKEGFQNTKASDP YG + N RNDSPSP Y + S F K +5F S ++YSKVY KE +D	120
Sbjct	61	YGSRNAFNI <mark>RN</mark> DSPSPTNYNLKSEFTKSPSNKAFSFGISREAYSKVYIKENPLADA	116
Query	121	NVPGPGTYKSVPTFGNKGRKFTLKGKLKDLPSSVKVPGPGAYKDLTTLPRTGRNFYSK +VPGPG Y+ P G + K+TL+ K ++ + PGPG Y L G F SK	178
Sbjct	117	SVPGPGQYQIPPIVGKEALKYTLRPKTQNPYTQTYKGQPGPGQYDTKPALNDKGLYFNSK	176
Query	179	FKSVCNGAIGPPTNTRFKDQTQTSAEIPGPGQYTPKTGLTGTGHYFLS FK+ +I PP++ RFKD 0 +PGPG Y P + TG YF+S	226
Sbjct	177	FKNSSATSIDPPSSVRFKDINGKLVFDVNNIQILGKRSVPGPGTYQPNIEMNKTGSYFVS	236
Query	227	RFKSSGTTVLGGARRASTHRVSNDGPGPGSYILPSDFGY 265 F+SS T +LG + + PGPG+Y LPSDFGY	
Sbjct	237	NFQSSMCRSHYHFDRQTNILGSTMKGTPGPGNYRLPSDFGY 277	

Supplementary Note Figure 3. Supporting information for +1 frameshifting at AGG_TAA. See Supplementary Fig. S8 for the legend.



MSLPRY LNNQKC DHKVLY EYYIPG DQLIKI FLEKCL	VANAC VIIKI KQFKD KEYNV VEICO QYDRE	USPELIAITIAKIYA VOKPESYYDYTNEELSMOPMDYYEVVQKIGRGKYSEVEDGVNT LKPIKLEKMQREIKILQTLYGGKNIIKLYDMAQDDVSEVTALVERVNHT FDIRYYIYEVLLGLDYCHSLGIMHRDIXPHNIMIDHEQRQLRIIDMGLA KRVASRYYKGPELLVDDRLYMYSLDMHSLGCTMANMHFQNPMFKGVDND IDGLMDYLKTYKLEINRYHQKHLKNWEKVSWEEFITKKNKHLVTKEALD KRIMPQEAIEHEYFAPVIEYKKKIHGGGEEVKTEE	
>gb EJ Length	Y6941 =333	1.1 Casein kinase II subunit alpha [Oxytricha trifallax]	
Score Ident	= ities	392 bits (1008), Expect = 3e-132, Method: Compositional matri = 181/330 (55%), Positives = 249/330 (75%), Gaps = 2/330 (1%)	x <mark>adjust</mark> .
Query	1	MSLPRYYANACVDKPESYYDYTNFELSWGPMDYYEVVQKIGRGKYSEVFDGVNTLNNQKC M+IP+YYAN C + P Y DY N+F+ +G + YF+++KIGRGKYSEV++G+NTLNN++	60
Sbjct	1	MNLPKYYANVCEEMPPEYSDYENYEVKFGSQENYEIIKKIGRGKYSEVYEGINTLNNERI	60
Query	61	VIKILKPIKLEKMQREIKILQTLYGGKNIIKLYDMAQDDVSEVTALVFERVNHTDHKV VIKILKP+K K++REIKILQTL G NII L D+ +D +++ AL+ E V+ D +	118
Sbjct	61	VIKILKPVKKTKIRREIKILQTLKNGINIINLIDVVRDPMTKTPALIMEYVDTGDVDFRT	120
Query	119	LYKQFKDFDIRYYIYEVLLGLDYCHSLGIMHRDIKPHNIMIDHEQRQLRIIDWGLAEYYI LYK F DFDIRYY++E+L LD+CHS GI HRD+KPHNIMIDH R+LR+IDWGLAE+Y	178
Sbjct	121	LYKSFTDFDIRYYMFEILKALDFCHSKGITHRDVKPHNIMIDHASRKLRLIDWGLAEFYH	180
Query	179	PGKEYNVRVASRYYKGPELLVDDRLYNYSLDMWSLGCTMANMMFQRDPMFKGVDNDDQLI PG+EYNVRVASRY+KGPELL+D + Y+YSLD+WSLGC + M+FQ++P F+G DN DQL+	238
Sbjct	181	PGQEYNVRVASRYFKGPELLIDLQTYDYSLDIWSLGCMFSGMIFQKEPFFQGKDNYDQLV	240
Query	239	KIVEICGIDGLMDYLKTYKLEINRYHQKHLKNWEKVSWEEFITKKNKHLVTKEALDFLEK KI ++ G + L Y++ Y + ++ ++ L K W +FI N+HLV++EALD L K	298
Sbjct	241	KIAKVLGTEELYAYIEKYNVTLDSHYDDILGQHTKKPWHKFINSANEHLVSEEALDLLSK	300
Query	299	CLQYDREKRIMPQEAIEHEYFAPVIEYKKK 328 L+YD +RI+P++A++H YF PV E+ K	
Sbjct	301	MLKYDHAERIVPKDAMDHPYFKPVKEFHAK 330	

Supplementary Note Figure 4. Supporting information for +1 frameshifting at ATT_TAA. See Supplementary Fig. S8 for the legend.



>comp6317_c0_seq1_GAATAA(+2/-1) MSSYMKKSSELEELKIELNSLKVDEQKEAVKQVIAMMTIGKDVSGLFPHVTKCI LSPSIELKKLVVLVIINVAKSKPDLTLMAVSAFTKDAHEKSNPLIRALAVRTMCCTRIEI KIATVLCESLKOLVDDPVVKKTAAISVAKTYINPHFTKELGFIKLIQGLOBGMAIV VANAVAALFEISRVAGKNYLKANKETIGKLLNALNETNEWGQIYILESIINYKPKEQKEA EEIIERIMPRLQHANPAVVLGATKNVLHFLKFVNLKSNKTTILKKLSAPLITLLSSEPFI QYIALCHLLILUQPINVFEKNVKHFCRFSDPTVYKLAKLOVMVGVADNTVNDIITEL HEVCNNIDQDFVRSVKAICQVVVKVDRVAKKGVEALREHVNQEQGSDSALQEAVIVASK ILRKYPKKFGLVKDIVKQQBRIDEPESKSAFIWILGEYSKKIEDAGEKLQVVISFJFDE NINVCLQIITSAVKHFIKOSDNYEDHVMIVLLASESSANPDLDRRGVJYMRMLSTDPSQ TKDTVLAKRPEVEEDLTKLNDDETBDIFIDIFISDTKHSALRPEKTASAPEPDSDEEVEE EEKPKKSKKKOKSKKOKKVKEETEKEETLEEDEEVTEDEPKDDLDDIFGLGIGDDPSN TKDTVLAKRPEVEEDLTKLMDDETRDIFIDIFIDIFISDTKHSALRPEKTASAPEPDSDEEVEE EEKPKKSKKKOKKSKKQKKVKEETKEETELEPEVTEEDKFDDIDDIFGLGIGOPUSN DEPAVDPLAGILDEGNGGGESTQAASPWDDNGLFGGFGASGSEASLFIKSEHAEVLSSST PGSQNKAAGLQIKARFYREGTSIKLDMIFVNSTAGIISDFDIHIMKNPFGLKPOFISVIF ISAGOFFITVECSIDQSANDLKNPPGCPVVJQTAIKSLDVVFQVPCLHTLLQGTPV AVTQTQCQQMANSIANKHSFTVSSARFAGSASDLKTRMQSNSFYPIYDELNSQIFATSTV NNIPILLRCTPEGSDIQIIACTPVAPLYQLIEEAIKEVISK >emb|CDW87346.1| ap-2 complex subunit [Stylonychia lemnae] Length=1023 Score = 706 bits (1821), Expect = 0.0, Method: Compositional matrix adjust. Identities = 428/1024 (42%), Positives = 610/1024 (60%), Gaps = 93/1024 (9%) MSSYMK--KSSELEELKIELNSLKVDEQKEAVKQVIAMMTIGKDVSGLFPHVTKCILSPS 58 Ouery 1 + KSSEL EL+ ELNSLK +E++EA KOVIAMMTIGKDVS LFPH+ KC+ MTNYFQNMKSSELAELQHELNSLKPEEKREAAKQVIAMMTIGKDVSSLFPHMVKCMETTQ 60 Sbict 1 IELKKLVYLYIINYAKSKPDLTLMAVSAFTKDAHEKSNPLIRALAVRTMGCIRIEKIATY 118 Query 59 +ELKKLVYLYIINYAK KPDLT+MAV++F KD+ + +P++RALAVRTMGCIR+E+I MELKKLVYLYIINYAKVKPDLTIMAVNSFQKDSRDIQSPMMRALAVRTMGCIRVERITEY 120 Sbjct 61 LCESLKDCLVDDDPYVKKTAAISVAKIYHTMPEHTKELGFIKLLQGLLQDGNAIVVANAV 178 +CESLK+ L D DPYVKKTAA+ VAK++ T P K+ IK+LQG+L DGNA+VVANA MCESLKERLNDQDPYVKKTAALGVAKLFQTSPRLVKDHSLIKILQGMLYDGNAVVVANAA 180 Query 119 Sbjct 121 AALFEISRVAGKNYLK-ANKETIGKLLNALNETNEWGQIYILESIINYKPKEQKEAEEII 237 A+L EISR +GKNYL+ N + + KLL ALN+ NEWG+IYILE I +Y + KE+E I+ ASLLEISRASGKNYLRLKNDQGLNKLLIALNDANEWGKIYILEGISSYDTSDSKESENIV 240 Query 179 Sbjct 181 ERIMPRLQHANPAVVLGATKNVLHFLKFVNLKSNKTTILKKLSAPLITLLSSEPEIQYIA 297 Query 238 ER++P L H NPAV+L A K VL F+ V+ + I+KKL PLITLLS+E EIQY+A ERVLPMLTHNNPAVILSAVKTVLKFMNNVSTQDLLKGIIKKLGPPLITLLSTEAEIQYVA 300 Sbjct 241 LCNILLILQQIPNVFEKNVKMFFCRFSDPIYVKLAKLDVMVGVADNTNVDIIITELHEYC 357 Query 298 L NI ILQ+ ++FE+NV++FFC+++DP+YVKL K+D++V VAD+ NV+ I+ EL EY LRNINFILQKYSHLFEQNVRVFFCKYNDPVYVKLEKIDILVKVADDKNVETILAELKEYS 360 Sbjct 301 NNIDQDFVRRSVKAIGQVVVKVDRVAKKGVEALREHVNQEQGSDSALQEAVIVASKILRK 417 Ouerv 358 +ID + V++SV+AIGQ+++KVD+ A K VE + E V QG + +QEAVIVA I RK GDIDPELVKKSVRAIGQIILKVDKAASKAVEIIHEIVT--QGGEIGVQEAVIVAKDIFRK 418 Sbjct 361

Supplementary Note Figure 5. Supporting information for +2 frameshifting at GAA_TAA. See Supplementary Fig. S8 for the legend.



>COMP6 MYSTKF TFSNIR EIMTIV KCTALG HYEFYQ IEKADI GALEYD >emb[C	951_CI RRVMTI EEDGQU DANQDI LKQMSI FDVDEI IKMKI EVVDI DW759	9_seq1_GAGTAA(+2 MAPLLLANPALALCEEPSTADRIRGNYENKIRFFAAPEIKIFE VYMSYQDFFHSLTPYNFVASKDDDDDDDDEENKDKEKEPGYFDKFTP KKIDFNEYIFFITLLQLPEGEVMRIIEKVNPEERKINKAQFAKYLTKLR SSFMPDGRKISTDEDHISKTILHLFNDKEYITIEDFCELKSKLKHALL DETISAESFAKSLLSCLNYYQASKYSRRIHSLKLEGRVSFKEYVAFHNL STYRFLSLGMFRDLCDDFAKLDPYCNQNKVSISDTQIATFFKVLDEDEN LEGKKNIGLGKEDKFKREMMEKIDRYIKKFQKYVGWT 18.1] calciumbinding atopy-related autoantigen [Stylonychia l	emnae]
Length	=426		
Score Ident	= ities	311 bits (798), Expect = 2e-98, Method: Compositional matrix = 169/401 (42%), Positives = 246/401 (61%), Gaps = 27/401 (7%	adjust.)
Query	9	VMTMAPLLANPALALCEEPSTADRIRGNYENKIRFFAAPEKIFETFSNIREED M + PL L+ N L+ CEE DRIRGNYENKIRFF+ PEKIFETF++ + E	62
Sbjct	42	AMIITPLAFQMLILNNQHNLSQCEEAPRQDRIRGNYENKIRFFSPP <mark>EK</mark> IFETFASSKNEK	101
Query	63	GQVYMSYQDFFHSLTPYNFVASKDDDDDDDDEENKDKEKEKEPGYFDKFTPEIMTIVDAN	122
Sbjct	102	GDLVMSYSDFFRALTPYNHSEIKDSKP-YFDKYKPDILKVADSN	144
Query	123	QDKKIDFNEYIFFITLLQLPEGEVMRIIEKVNPEERKINKAQFAKYLTKLRKCTALGLKQ D I F E+ FFIT+LO+P G + + K E K+N+ +F+K LT LRK T LG KO	182
Sbjct	145	GDGVISFPEFFFFITILQMPLGLIHKEFTKHVKEGGKMNQDEFSKTLTTLRKKTLLGTKQ	204
Query	183	MSKSFMPDGRKISTDEDHISKTILLHLFNDKEYITIEDFCELKSKLKHALLHYEFYQ	239
Sbjct	205	INKGMVPDARLISATEDDFSQTNNEICNQLFKNKTLFSYEDFINFRDELKIALRHYEFHQ	264
Query	240	FDVDE-DETISAESFAKSLLSCLNYTQASKYSRRIHSLKLEGRVSFKEYVAFHNLIEKAD	298
Sbjct	265	++V+E +++IS E F KSL+ CL Y +A Y +R+H LKL+G VSFKE++AF I+ D YEVNEENDSISMEDFTKSLMVCLPYKEAHMYIKRVHELKLDGEVSFKEFLAFQRFIDDVD	324
Query	299	IIKMKISTYRFLSLGMFRDLCDDFAKLDPYCNQNKVSISDTQIATFFKVLDEDENGALEY IK K+ YR+++L + LC +F + D +C + V I+ 0+ ++LD D NG L++	358
Sbjct	325	HIKEKVLVYRYITLDQLKSLCKEFCEEDEFCKKENVQINPKQVEALVRLLDLDGNGQLDH	384
Query	359	DEVVDILEGKKNIGLGKEDKFKREMMEKIDRYIKKFQKYVG 399	
Sbjct	385	DEVIGVLDQRQLLGQGKENELKEAIESSFKKVVQWFRETLG 425	

Supplementary Note Figure 6. Supporting information for +2 frameshifting at GAG_TAA. See Supplementary Fig. S8 for the legend.



Supplementary Note Figure 7. Supporting information for +2 frameshifting at GTA_TAA. See Supplementary Fig. S8 for the legend.



DEQINKAENRAMEEKVYQILRDLENNEYSKVVQLEISLDEANEKYRRLENTINGCDVPLRKK INKLENQAEQTITWYQVYSKSVVLVDLEYKKKIKRKDEKTINGCDVPLRKK KKILRGLKSLKWKANDENRILEGNNVTGIPSSGRVVKPLRGGRKDKPARIMSSIQS kinesin motor catalytic domain protein [Tetrahymena thermophila SB210] sequence ID: gb[EAR95613.3]Length: 916Number of Matches: 1 Related Information Range I: 8 to 898GenPeptGraphics Next Match Previous Match Alignment statistics for match #1 Score Expect Method Identities Positives S41 bits(1393) 3e-174 Compositional matrix adjust. 358/939(38%) S38/939(57%) Query 30 GQGNICVVCRFRPLNQNELNHGCSNVCADFHPNKKSVTIVTEGDG---VTNKNEFTFDRV 86 GCNI VVCR RP N++EL GS C + F + T G NK F FDRV Sbjct 8 GSGNIQVVCRVRPFNKSEL-ENGSVPCVEFLDQQTIRVKLINTDGKEKADNKQLFNFDRV 66 Query 87 FDINSQTEVYNQAAKPIIESVMECFNGTVFAYQQTSGKTFTMQGDIEDIESQGIVPR 146 F++ ++Q ++Y AAKP++SVEGFNGTVFAYQQTSGKTFTMQGASIDEKLKGVIPR 126 Query 147 MVRTVFINIENSSENIEFTVKVSMMEIYMEKVRQLLDPTKANMKIKVDKHKSAVHDLTE 206 MV+TVF I ++ +HEF +K-S+HEIYMEKVRDLLDTKANKHKSVTIQUTE 186 Query 207 RYIGSDLDVVDIARIGNNRKVASTSMNDQSSRSHSIFVMTVHQNNLDDQTSKTGILYLV 266 +Y++ +DV+D+HEIKINSVEIVMEKTRDLLDTKNHMKJNDLAKATGKLILV 246 Query 267 RVIGSDLDVVDIARIGNNRKVASTSMNDGSSRSHSIFVMTVHQNNLDDATSKTGILYLV 266 +Y++ DV+D0+RIGN NR V +T+MH× SSRSH FFM+V QNNL+D ++KTG L LV 470 LAGSEKVAKTGASGHTLDEAKGINKSLSILGKVINALTDGKSKHIPYRESKLTRILSES 326 DLAGSEKV KTGA G LDEAK IN+SLSHG VINALTDGKSKHIPYRESKLTRILSES 326 Query 327 LGGNARTALIITCSPSVYNDMEISLEGFG AF INNK KVNKLETVAEMKKLLSSES 386 Query 327 LGGNARTALIITCSPSVYNDMEISLEGFG AF INNK KVNKLETVAEMKKLLSSES 386 Query 387 IIEVRMYRWLEGGTEILEGGEVFEDEYKDUGLSSKAPAPAKEEVKE-PAPSKEPEQDDD 444 +E + R MUE I + LG +PE + + + +EE + PA E * 5bjct 367 QLEEKTRRVAQLEDYIQQLGSQLPSETNLNQQEDQTQIIDSQEEIPQNPAQISNIEEIVR 426

Gaps 144/939(15%)

Supplementary Note Figure 8. Supporting information for +1 frameshifting at GTT_TAA. See Supplementary Fig. S8 for the legend.



Supplementary Note Figure 9. Supporting information for +2 frameshifting at TTA_TAA. See Supplementary Fig. S8 for the legend.



Supplementary Note Figure 10. Supporting information for +1 frameshifting at TTT_TAA. See Supplementary Fig. S8 for the legend.

SUPPLEMENTARY NOTE 4. IGV screenshots of ribo-seq reads

alignments in the vicinity of selected frameshifting sites





	comp5116_c0_seq1					
		1	140 bp	- 41 bp	150 bp	160 bp
remappingCrassusRibo, sorted b Coverage remappingCrassusRibo, sorted b			_			
Sequence A →			T N		A T G A A A T C G C N E K S R	
	comp5528_c0_seq1	140 bp		- 30 bp		150 bp
remappingCrassusRibo.sorted.b Coverage	[0-12]					
remappingCrassusRibo.sorted.b						
						-
Sequence A 🔿		F T A C A G		T A G A		

	comp5973_c0_seq1				
			32 bp		
	1.750 bp	1,760 bp	I	1,770 bp	
remappingCrassusRibo.sorted.b	[0 · 20]				
Coverage					
				1	
	(
					0
remappingcrassoshibb.sorted.b					
Sequence A =>	A A T G A G G T	GTACTT	A T A G G	C T T G A T	G G A G G A A
20		V Y L			E R E N
	19 N N		N R	L	6 6
	comp6034_c0_seq2				
	comp6034_c0_seq2		C		
	comp6034_c0_seq2 				
	Comp6034_c0_seq2	I			
remapingCressaRilos.sorted b	comp6034_c0_seq2		30 bp		500 bp
remappingCressusRike.sorted.b	comp6034_c0_seq2				530 bp
remappingCrassuffiles.sorted.b Coverage	Comp6034_c0_seq2		30 bp		539 by
remappingCrassuRibo.sorted.b Coverage	Comp6034_c0_seq2				530 by
remappingCressuRibo.sorted Coverage remappingCressuRibo.sorted	comp6034_c0_seq2		30 bp 580 bp 1		530 bp
remappingCrassuRibo.sorted b Coverage remappingCrassuRibo.sorted b	Comp6034_c0_seq2		30 bp		530 kp
remappingCrassuRibo.sorted.b Coverage remappingCrassuRibo.sorted.b	(0.37) (0.37)		30 bp		
remappingCrassusRibo.sorted.b Coverage remappingCrassusRibo.sorted.b	Comp6034_c0_seq2		30 bp		530 kp
remapingCrassaRibs sorted b Coverage remapingCrassaRibs sorted b	Comp0034_c0_seq2				530 kp
TermappingCrassusRiko sorted b Coverage remappingCrassusRiko sorted b	Comp6034_c0_seq2				530 bp
TermaphingCrassaRibo sorted b Coverage remaphingCrassaRibo sorted b	Comp6034_c0_seq2				530 bp
remappingCrassuRibo.sorted b Coverage remappingCrassuRibo.sorted b	comp6034_c0_seq2				530 kp
TemapingCrasuaRibo.sorted b Coverage remapingCrasuaRibo.sorted b	Comp6034_c0_seq2				500 bp
TemapingCrasuaRibo.sorted b Coverage remapingCrasuaRibo.sorted b	Comp6034_c0_seq2				500 bp
remappingCrassusRibo sorted b Coverage remappingCrassusRibo sorted b	Comp6034_c0_seq2				500 bp
remappingCrassusRibo sorted b Coverage remappingCrassusRibo sorted b	Comp6034_c0_seq2				500 bp
remappingCrassusRibo.sorted Coverage remappingCrassusRibo.sorted Sequence A →	Comp6034_c0_seq2		30 bp 500 bp 1		530 bp
remappingCrassusRibo.sorted b Coverage remappingCrassusRibo.sorted b	Comp6034_c0_seq2		30 bp 580 bp 1		530 bp
remappingCrassusRibo.sorted b Coverage remappingCrassusRibo.sorted b	Comp6034_c0_seq2		A T A A C		539 by 1 539 by 1 7 7 7 7 7 7 7 7 7 7 7 7 7
remappingCrassusRibo.sorted b Coverage remappingCrassusRibo.sorted b	Comp6034_c0_seq2				
TermappingCrassusRiko.sorted.b Coverage rermappingCrassusRiko.sorted.b	Comp0034_c0_seq2				
TermappingCrassusRibo.sorted b Coverage remappingCrassusRibo.sorted b	Comp6034_c0_scq2		30 bp		530 bp
TermaphingCrassusRiko.sorted b Coverage remaphingCrassusRiko.sorted b					530 bp
TemapingCrasusRike sorted b Coverage remapingCrasusRike sorted b					

	comp6054_c0_seq1			
			— 54 bp	
	1,250 bp	1,260 bp	1,270 bp 1,280 bp	3,290 bp
	(0 - 29)			
Coverage				
	· · · · · · · · · · · · · · · · · · ·			· · · · · · · · · · · · · · · · · · ·
		-0		
remappingCrassusRibo.sorted.b				
13 184				
Sequence A =>	AAGTTTTGATAGA	AGAGAGAAAAA	A T A A G C A A G A G G A A A	A T A A C T T G A A A C G G
	K V L I R	R E E K K K		
	the second by the second se			
	comp6951_c0_seq1			*
	comp6951_c0_seq1			A
	comp6951_c0_seq1_	160 bp	- 41 bp	130 bp
	comp5551_C0_seq1	160 bp	- 41 bp	189 bp
remappingCrassusRibo.sorted.b Coverage	comp5551_c0_seq1	160 bp	- 41 bp	130 bp
remappingCrassusRibo, sorted b Coverage	comp5551_c0_seq1	360 bp	- 41 bp	
remappingCrassuRiba.sorted.b Coverage	comp5551_c0_seq1	360 bp	- 41 bp	
remappingCrassusRibo.sorted. Coverage	comp5551_c0_seq1	360 bp	- 41 bp	
remappingCrassusRibo.sorted. Coverage	comp5551_c0_seq1	360 bp	- 41 bp	
remappingCrasusRibo.sorted.b Coverage	Comp5551_C0_seq1	160 bp	- 41 bp	
remappingCrasuaRibo sorted Coverage	comp5551_c0_seq1	100 bp	- 41 bp	
remappingCrassuRiba sorted b Coverage remappingCrassuRiba sorted b	comp5551_c0_seq1	360 bp	- 41 bp	
remappingCrassuRibo.sorted b Coverage remappingCrassuRibo.sorted b	comp5551_00_seq1		- 41 bp	
remappingCrassusRibo.sorted b Coverage remappingCrassusRibo.sorted b	comp5551_c0_seq1		- 41 bp	
remappingCrassusRibo.sorted.b Coverage remappingCrassusRibo.sorted.b	comp5551_c0_seq1		- 41 bp	
remappingCrassufilibo.sorted b Coverage remappingCrassufilibo.sorted b	comp5551_c0_seq1		- 41 bp	
remappingCrassusRibo.sorted.b	comp5551_c0_seq1	160 bp	- 41 bp	
remappingCrassuRiba.sortedb Coverage remappingCrassuRiba.sortedb	comp5551_c0_seq1	160 bp	- 41 bp	
remappingCrasusRiba.sorted b Coverage remappingCrasusRiba.sorted b	comp5551_c0_seq1	100 bp	- 41 bp	
remapingCrasusRiba.sotedb Coverage remapingCrasusRiba.sotedb Sequence A →				
remappingCrassuRiba.soted b Coverage remappingCrassuRiba.soted b Sequence A →	Comp5551_C0_seq1			
remappingCrassuRibo.sorted b Coverage remappingCrassuRibo.sorted b	Comp5551_C0_seq1		- 41 bp 10 kp 1 kpp 1 kpp	
remappingCrassuRibo.sortedt Coverage remappingCrassuRibo.sortedt	Comp5551_00 seq1		- 41 bp	
remappingCrassuRibo.sortedb Coverage remappingCrassuRibo.sortedb Sequence A	Comp5551_C0_seq1.		G T A A A A A T C T F	
remappingCrassuRibo sorted b Coverage remappingCrassuRibo sorted b			- 41 bp	
remappingCrassuRibo sorted b Coverage remappingCrassuRibo sorted b	Comp5551_C0_seq1		- 41 bp	
remappingCrassufilito.sortedb Coverage remappingCrassufilito.sortedb	Comp5551_C0_seq1		- 41 bp 10% bp 1 b 1 b 1 b 1 b 1 b 1 b 1 b 1 b	
remappingCrassusRibo.sortedb Coverage remappingCrassusRibo.sortedb	Comp5551_C0_seq1		- 41 bp 10% bp 1 b 1 b 1 b 1 b 1 b 1 b 1 b 1 b	

mp7073_c0_seq1 — 37 bp 1,690 bp 1,700 bp 1,710 bp remappingCr Coverage comp7194_c0_seq1 — 52 bp 230 bp 210 bp 240 bp 250 bp 220 bp [0 - 102] mappir overage = = .



comp7880_c0_seq1 46 bp **1,190 bp** | 1,160 bp 1,170 bp 1,180 bp remapping0 Coverage [0 - 29] comp7882_c0_seq1 47 bp 1,190 bp 1,210 bp 1,230 bp 1,220 bp 1,200 bp [0 - 89] emapping overage = 11

