

Online Supporting Information A

The benchmark data for predicting the cleavage sites in proteins by HIV-1 and HIV-2 proteases (from K.C. Chou, Prediction of HIV protease cleavage sites in proteins, [Analytical Biochemistry](#), 1996, 233, 1-14)

(1A) \S_1^+ : list of 64 cleavable octapeptides by HIV-1 protease

CANLSTFA
VVIATVIV
TQIMFETF
GQVNYYEEF
PFIFEEEP
SFNFPQIT
DTVLEEMS
ARVLAEAM
AEELAEIF
SLNLRETN
ATIMMQRG
AECFRIFD
DQILIEIC
DDLFFEAD
YEEFVQMM
PIVGAETF
TLNFPISP
REAFRVFD
AETFYVDK
AQTFYVNL
PTLLTEAP
SFIGMESA
DAINTEFK
QITLWQRP
ELEFPEGG
ANLAEEA
SQNYPIVQ
PGNFLQSR
KLVFFAE
GDALLERN
KELYPLTS
RQANFLGK
SRSLYASS
AEAMSQVT
RKILFLDG
GSHLVEAL
GGVYATRS
FRSGVETT
VEVAEEEE
LPVNGEFS
ETTALVCD
HLVEALYL
HYGFPTYG
DSADAED
GWILGEHG

GWILAEHG
QAIYLALQ
EKVYLAJV
VEI CTEME
TQDFWEVQ
LWMGYELH
GDAYFSVP
ELELAENR
SKDLIAEI
LEVNI VTD
GGNYPVQH
ARLMAEAL
PFAAAQQR
PRNFPVAQ
GLAAPQFS
SLNLPVAK
AETFYTDG
RQVLFLEK
QMIFE EH G

(1B) S₁⁻ : list of 239 non-cleavable octapeptides by HIV-1 protease

KVFGRC EL
VFGRC E LA
FGRCE LAA
GRCE LAAA
RCE LAAAM
CE LAAAM K
ELAAAM KR
LAAAM KR H
AAAM KR HG
AAM KR HG L
AM KR HG LD
MK RH GL DN
KR HG LD NY
RH GL D NY R
H GL D NY RG
GL D NY RG Y
LD NY RG Y S
D NY RG Y S L
N Y RG Y S L G
Y RG Y S L G N
R G Y S L G N W
G Y S L G N W V
Y S L G N W V C
S L G N W V C A
L G N W V C A A
G N W V C A A K
N W V C A A K F
W V C A A K F E
V C A A K F E S
C A A K F E S N
A A K F E S N F

AKFESNFN
KFESNFNT
FESNFNTQ
ESNFNTQA
SNFNTQAT
NFNTQATN
FNTQATNR
NTQATNRN
TQATNRNT
QATNRNTD
ATNRNTDG
TNRNTDGS
NRNTDGST
RNTDGSTD
NTDGSTDY
TDGSTDYG
DGSTDYGI
GSTDYGIL
STDYGILQ
TDYGILQI
DYGILQIN
YGILQINS
GILQINSR
ILQINSRW
LQINSRWW
QINSRWWC
INSRWWCN
NSRWWCND
SRWWCNDG
RWWCNDGR
WWCNDGRT
WCNDGRTP
CNDGRTPG
NDGRTPGS
DGRTPGSR
GRTPGSRN
RTPGSRNL
TPGSRNLC
PGSRNLGN
GSRNLGN
SRNLGNIP
RNLCNIPC
NLCNIPCS
LCNIPCSA
CNIPCSAL
NIPCSALL
IPCSALLS
PCSALLSS
CSALLSSD
SALLSSDI
ALLSSDIT
LLSSDITA
LSSDITAS
SSDITASV

SDITASVN
DITASVNC
ITASVNCA
TASVNCACK
ASVNCACK
SVNCACKI
VNCACKIV
NCAKKIVS
CAKKIVSD
AKKIVSDG
KKIVSDGN
KIVSDGNG
IVSDGNGM
VSDGNGMN
SDGNGMNA
DNGGMNAW
GNGMNAWV
NGMNAWVA
GMNAWVAW
MNAWVAWR
NAWVAWRN
AWVAWRNR
WVAWRNRC
VAWRNRCK
AWRNRCKG
WRNRCKGT
RNRCKGTD
NRCKGTDV
RCKGTDVQ
CKGTDVQA
KGTDVQAW
GTDVQAWI
TDVQAWIR
DVQAWIRG
VQAWIRGC
QAWIRGCR
AWIRGCRL
KETAAAKF
ETAAAKFE
TAAAKFER
AAAKFERQ
AAKFERQH
AKFERQHM
KFERQHMD
FERQHMDS
ERQHMDSS
RQHMDSST
QHMDSTS
HMDSSTA
MDSSTSAA
DSSTSAAS
SSTSAASS
STSAASSS
TSAASSSN

SAASSSNY
AASSSNYC
ASSSNYCN
SSSNYCNQ
SSNYCNQM
SNYCNQMM
NYCNQMMK
YCNQMMKS
CNQMMKSR
NQMMKSRN
QMMKSRNL
MMKSRNL
MKSRLTKD
KSRLTKDR
SRNLTKDR
RNLTDR
NLTKDRCK
LTKDRCKP
TKDRCKPV
KDRCKPVN
DRCKPVNT
RCKPVNTF
CKPVNTFV
KPVNTFVH
PVNTFVHE
VNTFVHES
NTFVHESL
TFVHESLA
FVHESLAD
VHESLADV
HESLADVQ
ESLADVQA
SLADVQAV
LADVQAVC
ADVQAVCS
DVQAVCSQ
VQAVCSQK
QAVCSQKN
AVCSQKNV
VCSQKNVA
CSQKNVAC
SQKNVACK
QKNVACKN
KNVACKNG
NVACKNGQ
VACKNGQT
ACKNGQTN
CKNGQTN
KNGQTN
NGQTN
GQTN
QTNCYQS
QTNCYQSY
TNCYQSYS
NCYQSYST

CYQSYSTM
YQSYSTMS
QSYSTMSI
SYSTMSIT
YSTMSITD
STMSITDC
TMSITDCR
MSITDCRE
SITDCRET
ITDCRETG
TDCRETGS
DCRETGSS
CRETGSSK
RETGSSKY
ETGSSKYP
TGSSKYPN
GSSKYPNC
SSKYPNCA
SKYPNCAY
KYPNCAYK
YPNCAYKT
PNCAKTT
NCAYKTTQ
CAYKTTQA
AYKTTQAN
YKTTQANK
KTTQANKH
TTQANKHI
TQANKHII
QANKHIIV
ANKHIIVA
NKHIIIVAC
KHIIVACE
HIIIVACEG
IIVACEGN
IVACEGNP
VACEGNPY
ACEGNPYV
CEGNPYVP
EGNPYVPV
GNPYVPVH
NPYVPVHF
PYVPVHFD
YVPVHFDA
VPVHFDAS
PVHFDASV

(2A) \mathbb{S}_2^+ : list of 22 cleavable octapeptides by HIV-2 protease

SQNYPIVQ
EEELAECF
TQIMFETP
GQVNYPEE

GGNYPVQH
PRNFPVAQ
AEELAEIF
PFAAAQQR
RQVLFLEK
ATIMMQRG
SLNLPVAK
ANLAEEA
PTLLTEAP
SFIGMESA
YEEFVQMM
RHVMTNLG
YISAAELR
GLAAPQFS
DGMGTIDF
GDALLERN
NPTEAELQ
RQAGFLGL

(2B) \mathbb{S}_2^- : list of 127 non-cleavable octapeptides by HIV-2 protease

KVFGRCHEL
VFGRCELA
FGRCELAA
GRCELAAA
RCELAAAM
CELAAMAK
ELAAAMKR
LAAAMKRH
AAAMKRHG
AAMKRHGL
AMKRHGLD
MKRHGLDN
KRHGLDNY
RHGLDNYR
HGLDNYRG
GLDNYRGY
LDNYRGYS
DNYRGYSL
NYRGYSLG
YRGYSLGN
RGYSLGNW
GYSLGNWV
YSLGNWVC
SLGNWVCA
LGNWVCAA
GNWVCAAK
NWVCAAKF
WVCAAKFE
VCAAKFES
CAAKFESN
AAKFESNF
AKFESNFN

KFESNFNT
FESNFNTQ
ESNFNTQA
SNFNTQAT
NFNTQATN
FNTQATNR
NTQATNRN
TQATNRNT
QATNRNTD
ATNRNTDG
TNRNTDGS
NRNTDGST
RNTDGSTD
NTDGSTDY
TDGSTDYG
DGSTDYGI
GSTDYGIL
STDYGILQ
TDYGILQI
DYGILQIN
YGILQINS
GILQINSR
ILQINSRW
LQINSRWW
QINSRWWC
INSRWWCN
NSRWWCND
SRWWCNDG
RWWCNDGR
WWCNDGRT
WCNDGRTP
CNDGRTPG
NDGRTPGS
DGRTPGSR
GRTPGSRN
RTPGSRNL
TPGSRNLC
PGSRNLCN
GSRNLNCNI
SRNLCNIP
RNLCNIPC
NLCNIPCS
LCNIPCSA
CNIPCSAL
NIPCSALL
IPCSALLS
PCSALLSS
CSALLSSD
SALLSSDI
ALLSSDIT
LLSSDITA
LSSDITAS
SSDITASV
SDITASVN

DITASVNC
ITASVNCA
TASVNCAK
ASVNCACK
SVNCACKI
VNCACKIV
NCAKKIVS
CAKKIVSD
AKKIVSDG
KKIVSDGN
KIVSDGNG
IVSDGNGM
VSDGNGMN
SDGNGMNA
DGMGMNAW
GNGMNAWV
NGMNAWVA
GMNAWVAW
MNAWVAWR
NAWVAWRN
AWVAWRNR
WVAWRNRC
VAWRNRCK
AWRNRCKG
WRNRCKGT
RNRCKGTD
NRCKGTDV
RCKGTDVQ
CKGTDVQA
KGTDVQAW
GTDVQAWI
TDVQAWIR
DVQAWIRG
VQAWIRGC
QAWIRGCR
AWIRGCRL
SQNYYIVQ
SQNYFIVQ
SQNYLIVQ
SQNYMIVQ
SQNYVIVQ
