# An improved deep learning method for predicting DNA-binding proteins based on contextual features in amino acid sequences 🔗

Nov 15, 2019

📖 PLOS One

Ruixiong Ma[1]

[1]USTB

1 | *Works for me*    dx.doi.org/10.17504/protocols.io.2rdgd26

Ruixiong Ma
USTB

ABSTRACT

With the explosively increased amount of newly discovered proteins, predicting the function of these proteins from amino acid sequences is becoming one of the main challenges in functional annotation of genomes. Nowadays a number of computational approaches have been developed to predict DNA-binding proteins effectively and accurately from amino acid sequences, such as SVM, DNABP and CNN-RNN. However, these methods do not consider the context in amino acid sequences, which makes it difficult for them to capture sequence features adequately. In this paper, we propose CNN-BiLSTM, a new method for predicting DNA-binding proteins, elaborately reconciling convolution neural network and bi-directional long short-term memory recurrent neural network. CNN-BiLSTM can explore the potential contextual relationships of amino acid sequences to obtain more features than traditional models. The experimental results show that the predication accuracy of the proposed CNN-BiLSTM method on the test set is 96.5%, which is 7.8% higher than that of SVM, 9.6% higher than that of DNABP and 3.7% higher than that of CNN-RNN respectively. Being tested on 20, 000 independent samples provided by UniProt that weren't involved in model training, the accuracy of CNN-BiLSTM is 94.5%, which is 12% higher than that of SVM, 4.9% higher than that of DNABP and 4% higher than that of CNN-RNN respectively. The model training process is visualized and compared with that of CNN-RNN, and it is found that the training process of CNN-BiLSTM support better generalization from the training data set, which shows that CNN-BiLSTM has a wider range of adaptations to protein sequences. On the independent samples set, CNN-BiLSTM presents better credibility, for its predicted scores are closer to the labels of the samples than those of CNN-RNN. Therefore, the proposed CNN-BiLSTM is a more powerful method for identifying DNA-binding proteins.

EXTERNAL LINK

https://doi.org/10.1371/journal.pone.0225317

📄 PDF

model structure.pdf

GUIDELINES

This is a method of recognizing DNA binding proteins by deep learning.

MATERIALS TEXT

WA+YsHNV4TofdDdin9k/OdbiNG3EBVYxFwYq0nBXvwhBgN/nEhso4Ps39rYkUa0htUN+Ae8FvecLVS9H474E2TNZWvp7D/8k7sT80QqTl
pk9QFDiNz+Lj5olg1DRwiF9YbSKaVeRV0gz7ejEd1cZK55ULzmPB/IuPFQDw2j60G4ecXKq+SGABNRW952c29b1708Cj/wM85S0yTzIdVy
2G2cYXyfG+WJTPWn1w9Ginw9Xn857bP04r2PpvXeW8WXjYCI0Xei651uSBnQbBGpCEtpHtG+haWlZ2aOu9c7Rjr2//dJnznNwZir35Wnt
48PZd+HpqxDTKOSE5TrLPmuHcinvNYZQVO5sV9HfBEeGq/hV8At/9XqympVN0ilxMsplS84u3JaYwhsZs+FrcVChubnXlSW9W8eU15A57
vt/s8Vgts+TnJgH+aW05NdBQW7odBKns4MMw2Ac6NqgClDYvQ1G3ZBIX2ajq+jLspz+mEuuflh+1A9cjzDT3+MIZcQoSNjAHEqgAApuWh
hdxUbKdIZkUwQh2Tlp/nJDnFLa5JGJ96M5TiWIycmcQuDCLoTqqHJnlE1PiL78ezc7dpFfwFz9HiuLLO8N0Lb/Nb/9tAyuUerSvuqU90Pp8
Wg80l8oucrn16/zTniloi+QCEe7bp64bG0QSaKtra6OoGfqYSAwEoemf8Z4bnu4M+qmdjivgwgUrwIZkiL0Sh38VN1YKgHIIo/zcilPfDVNzum
qvS3YwRDZ+qV5W7l2QZ61osIhlkKCdc/1/cCCO9MkSgBhOy5Xfwa62ty8uGSfD23HX3AgkcfBg2mfujyoInhus8IY3Wt26tIrA/RcArBpeXO
CmprtBF1+9olkMjScgTk7hcL9Nev1DmYgpJVvGThSsHomPZtHeUH5+0ADEPqqswnm55Ytf4UYdzxraQxqmWQ4dx0p/KigwfpCW0BPDW
NrNbgkV9lF+zpxLAUF+fh8bmDvb1EUfEXQbjkNAJqWG6ib4/d8tkeE6AIIfXDUReeWdU+tm7guTBvT6khsK66xQ8KUN5kFkCo11Kk8BrxZk
00PsSvoS4JY4MOnm1hDmaiXAxI9UGb83w7K0Zl9fwdfGiszXHTGyHAFlCYumess8+A1/cu8shXs2CkuXiIlFYHBXdTIIBAONP/Cwx64rlnh
qMtFr4ntSpG9Qd8JUwB7nC9VpqdUdDP+gr9B/+L4COtOIx51lqlcpXvaar7whTH8XWIr2zc9dPu9ioj25mzLELaWltdtXwO7xWzu6xsVVyz

n3FGJvp4Sn/BI/Id8qWGGnwQXJbEJ/KSmlEmN9o5lHV6kJ3cVJbt6Ld59G4uX4nQkrWgawSjdoWrojEhA9AJEuiKZLAhICMwOM8hYaP1
gFAirw01St3tOLWH5SHr3MWYRCJ35s2NqA+uchKIJKKhoP9Zt1QFXHdwOpkpPZdPjU01lJ1/gREWVcLq1kzM94G2jX7ZSCHB2FjYnJQN
Q+GUNPZYkkc1RXx6A44Gp/kDlqTN6hlppq1il57LdeX011k3C0A3ufQwt5vd/QmxBDqUnmONLh1HwCYd6JGrMn/bxKkRzmyO3GD946S
V1YMVI8ZF4lfZ5e76UMV7nIqJS6vnRtBkm+fkhonGazotYE9iqpuzTJSajDsKba/Dm4FWHSEmz+2JSTFw9Y6JZ9pezTq23cflaOn4d2dWY
AUiwUtd9JrG1wTyg9KqdX2WrRJrJmGgo5TjU1mVMd1sbqqdYbJnc3zUnyQ1npaatS4aFFHTF0HIHEy5VOWNU3yRhnFJw2UUTEK7qnh
eNGFmhdpPqmEuoHZLShiZ9VooHLN9hJH52+QD0fWtKkwAI8VZxsafh1WaK+z7m8eZgBrkxeHwstMHUbVWoiewQ/Ri7Z6KY8o/GFZZEk
hl7JL0i2BigsNI2WYZ22SC2FerilCVihT5+KhkiZaubnwF62EmDWNn8mTO7vWqgEQptn9hs+ybFdmnr+UcihUhb7EcIqxePAPN5859Jih3Ef
M/c9bqDKqKW1xCMT6RorsY8XGq9slTIjR/bd/kGCwB8J/2+9Wyk1I0O4ussoEURJ1pcG2r32KoqSh6hGPLSGh+wAxMlS7P2ILksD7G4Xq
zcYWi/nO8LY0/cg9k/rjq80zvU3Ak1vat6V45qfOq7QMZo+gDsVb/z8UKD4HJ3dHLm41g2i2+0Wfe7qfJ0s73OJv4fvhaJM8HHrL6w8zk0f
obki2UMs7dg595tbKvMe/F62NYN5sGfuSyqEFa892brfse/td+UYwO7HR7McFRxk33e1qVJgQ93ZiXwH7ompumS/kScCLs8H8GhWyZDm
BZjgp9sRh95TS72RLqH+bD6aF4kG2lwCWS46M5IunCZ+9ynZjrLv+85Erfl22SL3iUPNwg8J1s5QdVo44Qa/hsp9QtXQ9HRX79D2mJLLpl
8A/KZSsD9JxByak6WNHtpYRbp8VjA8thdC9roperJNUDfId8djKGfGbP4KteFBkie2taICjwLjv+kXd14hMEh7mHa1h5c8PRJCvetf+erRiRq6
UDM93BxRpKEWdExhQeEAK94bXWcKCmOhkS3BCst8mrFs77uO5xQK+PkTOgK2P0RURyTaT4cDfGSL/k4O/pnUoitEZk5nbtK2qk+5NVz
P+7A88Rj1RzXBdIO9S/PiR9Bot3V+xGzc5CU8dbfq9uaxiLim/zzmbAF+f6dlvSrEiAWJOpokiS3bCOaDHh558vqq6HWenuFIb5Vh8Qp1f4g
mL9ht/wdKbT1uMf+Jdxk+bO/1pwwNgx9Bn7RbS17LmZG7v8eBjgYYZtBXm1sBVNAArhNaXYJ/Eu9ke+1Za65+zdAwoYClA0ez8x4fPwL
fjuTSomy6vwTHt9yoQLMW+UF4uVtPxfZO+svYks3uZbyTEOmmu6GUEFCHBJZgExMHC14YHWelyJbnZgCzgICJLCUj1mRwWvBz2fs04A
Tm9WWWbcI1YeOfhrjXBO2sJrm3a6Jott8vlHghZDEbHoWHhMt1nZNSmPKHUb6NjuGZllCxyMqOM1ESSD1K7n9qYGP0E+Q9uwTLXOD
AnCaZUQj3uqsAqJmp9O/0ZrX4vV4ed4dMSS/KLTc6tKHla8nA9aCu3kqHGZDIdN+61RNWFdQ68ONpkCZeL8ifJQnl1sO2e6WBkpX9xdm
Ma9/cWuCHBVQw8VjicX75GrxZwZiux6E83eeeBNhG44Hn3gxyOooQZPwubVG+3suF4ok6pTjKN4YVcZ3Vtq45SbcSuxbNPKKKSZPtjLqY
GqzDmYlkEaxpsbyBbPvwCPQEVLW9rzO09+++wh9CyZz3JH5TaK1GTw94owgye/aVa+A3mIeoN2QxeOKAlFTeHCHQl+olDy0+GuZFzBO
M0viAkyA9UQgnz2KOyCmebqmcLdGsKsbdF8mSQIXF1B938A/pzUlfhDsIEzrbteBwM4kWC2jnGjtHtWmaiH1DxZXOYpocUrl5CvJau6WjA
ARIa/teCdklnUts1HOqNFihoeTnqVgVDW8YBSS0tPoplp69AQoHuFJ4d5P/JNcKEK2zHx5pSaGPS3vKJmNvPcMDSlD2JsEII6uC+82yaVn
6b2PaLZwdYNlYrD/nR9yD66ulgyI7t7Ut/65elT/zJCQ7bHrA4JdqfwS+/AU1R44sftUzRBET7GLjGj/ASJ4d6UWenl6PL/4c71AlW9VJz+G9X
TAmqaKxEl2/r++K1m0uaKfmHblaHPpDtFFTLAeb1u2fHtrJdtOt4xMtsl196UOD7pKMll6oCVYLHL1e5oo/kNpc6DzBQ7fn3tpPqcCZx6ZP8l
b5OOVqLtGkJ7YL1MxEGu+VhJSoAf7FzSlbQXKCWTr4kQhS5RyQvD/vpdE+UNBammCqtF1lRdeKHoduqLhAlZ+fQnb9q5L+4z0lTfwy8vq
FBN5mR2kRxc6avHOMnpiLsRjqH7T+0nngB9JXipF2VVysPfeOnBNGGEVS2xftbW7Y9dUUk227FLco4f/0qClP4twZhoVrPS1apH27VH1q
UnH/P2YqE2CzWIQrNETWNBL36DSxPlbvYxlO7k4ZV3U98CYALpNILlXcWYNlvBGn/Dg2DAqYehZKLBB9xpYMpo0RUKhXzjeB07p7h9sM
FAps38rlkF5DCSndAOKLcLSGo7nzIIlfC+FZqonh/ZVUoPDG6X3MhvvRXI5RiCB+F4OeU/Cr9/scBkX6/STVFYE4IE47NVctx9X3UPCu98D7
N2mX6hgnvDFaGnz5bQi6DAeKw6vRQARXSGUHlwsMGM2Y5yMtKysaA55et2vFh2nNzsV78IK9WxPKOUnBNIMsgDaiF0BsahxtfuYpBW
hjSHHbeaG2470RAbGsfWz25MkDadG/+xNLPyTvQTLi4qLY3CZWEw44LcF8r5UEFwdyvgJYvAJQzfgPnYjrGZJcUr/3bkyu9j52bXbMOeZzF
rHMdEcjbcgy7v16A/SDtia6+Ivgf+pXSIQDJ1LBRoH9rogRys8tOrFstFKZ+wKWkwqVl81bXvo1daxetQyY0jg2L2m/4SgSlVJGk5DpSUtdB4
8yO6Hxeis5tPUWPHz+Xa7SMS+W1Qzct+VpZly4r3fAOyuTBIoQ8GSdqdxTKwoEZxIdHquilfupXvFAMywdLsJDVil1jwhMpkAU7FBGO0WT
jm6A9NZIuJxgrCGFlCFO0/rv0TDyyfAX8N9+mpBDMHk3YiSASlWYgIzzEIfx/QCbsoaPQU/u80a+eEyYs0DT/c4TPPY/ge4cP84nPXLPyhJB
l+JH1a9gmLKAtlVBrc9m14XmRUsPBBXUWKse00F0IjYdTkoqS0O73WJlBmTF3cy8KJTd0LnJkZNZfFyp2eB1zKX0A1KQ7W4b0ck6QNGB
e+vazTOSrAmV8uOVa7RBIXQnhZwqQMud+sFP2tV1fpyUeogkuq9owiP8CZ0sj2OPz9MuNWbwJGi39rTRaYifPW8DVFQhc9Zml+iAxHHO
LoZkQe1DkvEJRV2zevQzWb0KsxtwSMcmSeZ3rDMczTLbazFrad7ErEf8iMXUgivCS4cC54PvKenJLoDBlfb9VE8r33JAt6RqpPSv+ptl424B
RyFbNm7JL9/4zKzsBN57y7KmVfWJJS5SGgdX4cTIwmpRLfuM3j3NBtSl0VUBL3UGSDpqRvMOLqlE6nJimrWHeGljTb7dGuOn8YJsW2Nif
sertdKb+NvAf3PHHAxywlBCgcsEdRjrV4+ysJbfnb8OK41vflN0dHO6Dnc1uhopdYXiOBJydLqOCcC7vkeu2d3z+kK4aRCrMucMFEhUnZ9+
1vamZKbyMPaxml4pCR9XI8JmchqG8rGHtyhU6mZ7RyZ7wEsNiCxl/9dLCQjd6P//+Q8YfBaY3thkfwJjydA6vluutxOOeTEQaqqxGKwgxq
mgfNx1BhvQHMWg+R+/8BeSx+dRfYMGknNgZ6lr1kCRVap+mlfwaiOqdcobrgHlTPCbhMzkTjevtSnaVG8old/idoRf34t5ZSlD68b/c4k0Sot
G3B/NwqYFHbhjOfWeUsiyrl0OKtpbjKWGq2X0io/FUbeY4uf1I4EzNB+tPJhVlo8m834S2DwI7QWB00WA/3Q3bofYXP55ad3wulvsNrKd1u
WBWhSUCEiVaEVwW7vdJ+PCGS4iCLbgF4nTaGRiSHiHSu3CltmgzUN++eftooa61JzlqUw/zKYxqpdGUw/njB0yT5AtTDsqIC4chBL6J8k8
6N7lzGr8Fbo6iauEXVufcGjR5TbPqoJP/m7VfHK7/8ml4cJwFHK4BmDJG65+brLjQssZGK9Wbp0pMDj8MAj5zZBJxV1tBu6f9J62Qp8O2
FoLCdCk7ApotuRMlK277CNUuifwbtQvhP3u91gqf2j58LeAAKOc3K8fABwjNijPZGtlAbvOK4kBOp/LMETGkx9KkTvxnQLO6KEnIfsRVuTZ
VUWkoXq1z+y2pg/wINIcnWNhCzL7aJM5NAJ86r+0cfaYofECs3GujdvoK/snxIl0uJdAGqDxEbKCft1Zfe1/WrFMEW7R0z0q46QJfqLnC5H
uCghA85zAVP0VMIEQlF/l19q9U9KyBvNNT0P7JnPfiHYJ+cpdC+k6nljjTyxfFdntRMm7d4VM2Q3TmGfzexXgYUF5YGaGj/eYHLJOwk2RL
DxuQl8NKAVs82Y/nCRQjtW21wXYORtQYDbOQk8ItVwrevpezD/QjiP+fkeXJD+3SM48YQRPodMZ1oypT+LnqqMx8ap4bhV74mSNJYlmB2
RxxbmwN7ZMQP5inssLv0AmtCXYh851OXOzWlx6dY52Yg1GYu8+U7syaoiWNYCAnWlpLHED/JI49RUrH0yq4UyrYoWlSp1dgKsGldRkP0
/kdilnkjUdXm3wXp0MZsPMsTnzKm65NHywpnZsOSlFV8ZS/HJA9pqv2WEVukI5ENG4Hku2AAt5RMK874FopoY336IJV+U44DOn14BHo
fmLB3oSMkE2i4MtdYTQfIGAa2UYST4I/VEPjxYCUcEichy7XJR08ijyhcDgNR/kwClmJ4Ea/QmPZBGKr2GDosXHpPQZgqGbh4KgUMIpv/j
WdXLOaNWh5yhedVRuvnWM1m8tgvTkcdxN17HRxKgDTvztQpPrirT7RTJLwejtGDjNwZa8H7ccC0ND8HnzkQkjzvkqU3Jors462/RDhrdfz
nfb8EKIPfNm9DywPC3cwxa/vaCGNRIlAr3iRYM7J1r6O/LjyMRpH8Jqn71V9XiFmZme+1IsknVrCYxKHyq9PGtjSODiYuXADl+TZv985JUC
ba5qQB69rD5ENpqxCobV63agSW714gVO8s3eED1v2k5UmoXvKReJ3ivOJQbfrUoolUKFA5uF0eO/tju61TwgkEqPiN3ENR7YiZE333sRKs
/RQyU4gjAaTffRlgML3dRMIjoQqtmPfN5iwvY+f3m9861BNLjfvackXyJgpLV792CNMT+C5j+UIJ5yMaRkfr5P2hvcgK8neXlqlGpgvQhNsAe
EbW1KQpJraalkodOmyf051x1UXDAXEtOiKuC/q2jNSCfxsBZhWl9qMai8wbOe1wioufEx+JSdDBDAVI/qZ3jlmQkDO7LTXIp3Bosz+hRBEx
636mHIC4+D60KRT9xzg4b401e5A8z47+0pu8vL0LK6FcOVIF0YgMlt5iCynHHZnPrDQ/XG9rj81zQCydObVBBqC0rvPZ1za5etG1TWPxri
BhhzO+efW+z8ne7kOMvXaqleHc1zySfaEj+rfQR8YGzjvuwp2+7XOIpE1TXcPtca1PQa91MeIOPtmg5WfhK2ayXOG6RlQ0FdVM8W6rrS9V
bZrl+EZC5XbvJangvvyqxa6rb9lUcO4TxTJoQgbv5SkzmzbOU5mTJUT06N9VnRkc19lBxcxMI4CtDUbUjoOZHlIMfM53KlKgU/1O8esObkY

X6r7hg6dJJyYSQfV8QYiJ01StTkmcBwdEwSxr7++gFDcSlyJtfYU4k1PxYCy8MRo+DSG2bbIFcQ9XECAtEFb1KJxRJbgQHJMtis5slfgOEbf
6ZTVWne7QYENXJwyYXB0B4mOhIYQOAsIYZum+91Imi32ES9+yjvk/rOeoMapk9mKB5N6x5L6s0uNc+tkuGGdwm3tz3K+qIu3f74xj1coD
uJPLQ1wK79PiaFkUpgk24XUolFGpnymIjkJ5qJEBy4zJ97AKlrG7uCy1YsJy7ovtKN5tuXPOi8E6JrCpB724bVFru6vMvxgJeWtWfaZ9LtQf
77rAnEOG9dVZI5PgkkCDYCj/kHJVSUsjhBXhdQkjXcEae4nV34j6KIthaZ8YbmHGKkO+ZMhIKXH4De5YDn/ONb7PjTEUkIH5pb7xQFPtorS
fw4KgmNn0srcBiPcbXjKXpKYGUBxp8vvQNosd7neHtlcrDhAruVw+iDvYEYXsOA1jojrlXGEkFtjuR219LWpMoio+b1Dv7w2vMlW2OekifE+g
oChLN3Kn6/zL8rdNQJRaueyEhmiXz0V8h4tWy7ZzvVdNUMMhx9e+kXRNDs6lmSGd9iqe5AGLD0t1EX61tQcJ/ikkaz0sZSZM0FdiNxABzz
rFXWy9BO3wGpV12wxhlTrxvxw9KZ6u97hoOV9fJlM8zMdC58LXll9/SpBMSzWCIYzH97MO2FzBXQ9+b/koxToTOauOBD87/7ZpnOoOp
Qv4XSVOrgdt2o7+oITxlcTJK36Echfknh85izCJGZ2D0gKTkgTlQaC2TkoEYh6OVZcla6kxGYYqk4vG/40zHWVAQNWzEhADnnhzxkQLPNZ
MAg0ro7p+JBc/5VoU3esl8Uul4xy8iuKpGduQXfENRKcAM+qFSIqWK51pK0XnprcRZnYIlowGyWBlWxsAolu8EX7mEFL8upZpdt3VD4Q2a
JLCeQWnLQtWqn8AzjKG/CX4BaECjO2gKVnqPhK6aOrQ6IE2XaW/6ypkB27PjGUsA1GIitc84Wk0syzD+oyCdfG0hJ3ObvR7Bie4+CMatO
4rg+LhFr1w8TN4bqM/Il/DRA9ClsYj+AlUAOADfN7A7b5OCrOaqU9pqqWLPwZkp2gFlfbeT5QUST4SzAcVqj9xrcofYfjogNhTa+55lEK8A97E
yy2Bsp46Aicv5zjHm799+PoZjkQf0UY3TFaotJQeAprWNJI8ifJEFwxlwelzaSfGdKOqX39AgL8PVLEb6AMBCAj0ul3cHNWKuf+hgzR5jU8l8
8D02ommpnGMm5Fp6dsg+O7IW4TJU3D4yawX9683LgWQnXhhYJ9VtUbWou5c14NUcHj9C6crxKC9vSdhiJWBhFdqz4jVpK3RwRwYm
wybfv713VjvTpO1ACb2TEvvUFOzmdKWK/wlQwaFY+TBibnNkTp4LEWt8IGHCtdzjnFPcty7aDi9DxBTNweS5snEh4u2aOWyIRexV29PUua
N5lgWKMt1oRZjyMG8xQgL5BzAZae9I2Iwfx4w2L6sUT1reWrI4X4JJbviIok7nvGcTOt45LlYRcr8Siq6O4gIVwGIHkAfJG768HXwyxbF8RFb
Fj9eG2fkxBe7MhPdX30cQa5kVWJMXcg1I1WOy1amOPYkSHhZqw7rGHrsU6cQiQKgOvsW6b3d2cako3CWDP3FjlpwxElNsmINhbd46op
RsKFgS1n5QQX0gc3XdU/fBLGXpJusbzgnNwo86XizH5QZBRlloFDJfuizc+ttBummVG+E6Y4M2glwoAAxUGH1yNH+Xtmrbh3aWlhERiR
Qt32w5yqow1N/S7n2o20iAYUOfDLD7O+42789PqQGqf3laM3LZvaPqAA82RB80mC9xYcLNGDGyj9tlB3U9XPjIsokS6Q64qs6ZAvHW/aR
73BeYNY4BdD2WBBobQMvMzCtiaLToknaJ4EiKiY1XfcsqJJixDDwkRUihtRDqXvmNjjJOPWjpbXmaZxVyBvnHy/jYyAbXpeMj277YcKX2iT
Gc54iisc1zZ464dwHBL1gbtxix1N31h5BMSBRrbZKXTGu7X3IqBX5eDzogOdj2Wqnju98qB5bQapBtbKSvv2VUWnY9YAVKuYhEWaeOCM
NWxY4BnXl0lKyIqMC2P6ZEulDELRcg2oWcVyLPAe124NjJC/GqYN2d0qHgFP9YDmZlOYOUs5PnQMQgkg03IoaX17ZB9SsD1ZgXnrRGb
Tk7DLiAFaAeQE8+qRc49SvczRO5TUjHZuV0H3WvnN8vQvcYJ1WqGzMyWQzE90AsYJNSAQbphKrSjBTZd6DkluDrWbnFhoT1rJ6oaQFds
xpAQSPjSDHOYVbmqQy0RlgpCZFtDFQ6+ryqGGsw4FsgUq9weF9DdVvMm1+dTA3EGXMm74EgPK0Eo8EuNRXrgBaREDwGnjGg7a90b
MSo/SQb5Z69zF7OIrFfYplkM+Lp9rlZc4wd2wvipwHoSIiOHXo/LKAj0xPecVyZRiSSR5dNUCMftq/3X9z9jTypOARtGKMAL6hMAddda8JP
h7G+/g7HPllf7QOWGPj6nafknKcECRa2H58BM/sTfXcXoiXrbc1gJFg63Ye78VqzOaSnmToOJ5oqHRMP236umIDbnN1n7ujuwfdQJvtTv
OYFce/nVp35goMk2udZ9/W4D25h8s0PewlM+qPrkHH9KSzayRP2AdWGyt+6XolLeySuOso5ktnJs/IKlksTPak4UnfBVskrk0A9rql5FwD9G
c+4Q0sHLPdHeZ00CJjUnNYtfQM8v11a0EiDQSRS23a0klDAZz7TYU2EcmEH/vTSkgO5XZmVHRrE9AuOEkulWoZZRFKjVshwdafXp021V
8CKAyLQMqfVH4r0Do5JXQXr3eBw7xHg00+j3TRbJJT6XXOZgX4P6KQZmGDnZKZtf4k6gOeVWel5OH4xV8B4t1of334ZJicVU73p+TNZ
xYcaCvhBaKXGCoTAnAphMK6E0IRQOQcEgLOBl0/Q9iB3M4SCtdP2rbBqdvPMbtvypXcAK3EmcRlQBiMpmQ2zb3wKkuhwpDlgw4GmQ
QzV8rbFo8jjg23wY32zOgpWLCS8ypOagbEIWR/WmQRooJ6Q2FoTNtx/f3iRoeELFTFuDCaZu49eK3QtldCeC7EAbGyTQU+UnGK8r1G+O
1g1vGy4Egs/nX3HcrkHz9oulTgk9nw+Rw5D9RN+8Vpl2UO3qp5RcpX+j/tZisDqCtupkLzBAIUsi1EGtEB9b+r/BZbjBFlw1t6oIizwH3BmdW/
qm00qLuWSwhKLZij80CsXtwI4yvVmT/h0KCl3Jn1djM81F4YvSz5eQxiulJnUYL43QTVZ1snKIplHt6I3RWhTRKBd/ztorAOXUZxz5kibvnTm
MY3J9tUs/UIP/oz9/z0RfCQ+3/7LI61Mk4CKMDjnGUTKFLyRvAwI0NlSeKGrRSqWzPV3adeQ//HjcBEY2T1P1wnsKQIBSxqT4Oa75WB8el
BmKJ7vxYlmS+OBjXREDG2TcwvYGYOWYNVUazksDp8JJOCBSR+LkkyVTCCkaDmVXsuKpq0RBS406oBlivhTGB0QOyznlzQOfxRVdVibG
8gxrYQ8ctf7c7w9wRhGnz14cU9u4Ew+SqirRn9MfGx9qrd4T+8IC2tzPaxvLZWAPk10FZ9qbO9/DEScFsOEsnHkc8wxzFAulQcvfo7qIbi9h
Y5JlIOgVaANaCRGlQ210lu+EsFr41/COAE2NKsMqewkszpbFg4hIvKV7VADRAuc77RfvLSyg0xGktOJV72aDOFvrxR+gVOtOkUNdmS8yD
maO3YLHf1D3ELZl6S0P9Q3Ou92MTfYi6ahfLn/TVUaCPBN0LQK7PxbDuuGKMwLsx9U8rS+zC2ugpLMksoCFShxAF9eMMc2+p5q94diy
Q0JGL6O+caC3/YG4nRcU9fIkaPApeb5Y2zIiAw/htYR+s+Bg17OARf2w52HDnztei6xVm1s6CbJba41C59rPgwBSr5KUEmj8OWIAiqWfMs
OBzWx4AP+yEqTveI7Sf625dweHIvFI/AFMczaPLYb566z6Eb5XNUCVEHrzD4Ok3we87uI6+y7qaSQbvo7wmTETm6mcm1eNQkPtNKbwb
t4nkMt+KrnR4wl5LxiY0xcnSdJR6uJkftyJu5ObIJYXj5EcrPZd3YncdKejIDNUCCiK/mRO6z4rA3sZvfqE39j7qlnrr8hqAcTlqo1ISS36gr/zph9
5grVbMpXS6lC/+ieUclb3/3giN2DDXOrWsEbi5kuNA12QfiYEjofdf0S4CZsOIQX4tdEm7dpqHPmmNJ/pgfzSmaCeFlpRxtRO4qMFOYdVCs
Ud6dIgFhcS2UYJ7C2uG6rfspDnFK+k0nY4cp1wG8PtfBIZ3Uuvc7rSn9YvbpyAZ6mjzomrHv90YsgY+HUmtn+vBIdRm+bBZWLqMOW78Ve
MMvw2fghqijuH8Cze3L1C6vsMIbddhqLJxIvcUT0q8Vzr5qovtWRdlHVGXK8IwitZ5j6sonnMoWFJFDxic5LQEGcJqgpU624vsPavxXPIONK
po1YXzrOaf9a+hZLvgDIMGmnsKc2EXMQYjqyAEtA5nasSwzoeGFQv/dWpjkEDsCl209117DbTkRpc3MGq7IBgjBaFFXYISMtsoQVPTEx2/
8EuQPWugVwCz1zHGJH/8jj7xE7Xd4XnPuDWyqp+rTNR21GFvWl1vNFuJR53j6UbzhJPux+EqPJIphJTsqxr03zv9NxMNVX6h3RxXEF6X
B51Qf5zRRQdCaVHllh+Z0IsKdEsFVDt8NjY0+4fKOTanF+qyiwwt9TWnimsqPRzW9rq+3lhQmTjAEidk4lJpR2IPOnsZXCf06+GEKoZ4LYLtZ
n4LZq6tVbs4erYHbo1cLHjBfHNdirGvNrYQ7/1TZPgqxUxNm8Xe3ahYpD8/Nud6IClQHOt27Kv7DD+rTTd6IuMeOKEn92P8ztQitw/9zjPijH
Co67yfzNPaqvNGL5L17myudNKuTkrT7GRKdoN01PcpWL2+PEqrDVjtvOeBZrykmkK76tsKV2sMai+NP1SmUuz998ztMxeHPsiN52XXEc
KDCE+xBsNuAqi/mrW3axZVKK7mccafCG3cJSDZV4L1+iC7TKp1gZrcjjj5ivZ9EjVzVapPT6SFuKcFz+SM94YKuMZVaj5CUVAS5EZ+pUn9c
DBWinr5DclvWRYMsUMaaJ1LFO/SqAVIUriXJM0cqVPm/n93hAOKoAuwN91QrtmOv6nCxDGjK6ABTc4TQEaPz7WLIN7CF8gMLJGaFlyu
N0/E2PUW9rh54doMzGFxwIDXclFxRXU+mszowYn5ECtLpMZ03WdCzCx4QoV6CuHqM7qxO7EmQosPf7Ahre0V0IvnrU22zIyfvQ8luBdm
hTeUmWIqGfgiEcQq7D3+bvIZuy4mjmv/YFC5RzG8cpuo0513r6qqhdJyxkD6l/7sFR2mc4btdYJQAz47zFV5z2KEvPlloi8pLW3Y/DmGYauh
sI3Hv9Y1ypPJNPyev6skRID0L0mbNzU8nYoW14vJnDohIdU2bOxtzIe7pnuFJYur+KzVXpsfooQ4vPHvD33Wa3Ff9+0dwrYLRxpQbziSs/6
MEkRCEmvvm481JgpUjEq1dKwcvR2jCJs5HlG3SIzW0+kJqCg0tkk2muWBBpG6BJNceO1mEBU+IB4sJN1BRJagnXkOgDneOMAX7UML
PipC5zpWbOpUPj4SnogHAsEMQsirjVT2nPndaYoWS2kVQXmpDdaVJ8Vok2YB7pq4EiButC3ec+OwSQwxPyFM9l9S/XJLOxbLYvNCYtqS
WQoG2GCLzbd2ZMRfhVxwLNqBPm7x2PPzKlGdNY/66eaCBVgWPot4/bLXCUL0HvNHFYtQrjklPE+wnnz7ueoL2I26DJzcHFTIJpsvODBg
qkZICWd2pOKbvWaBbRJVfNCOmNM9+Kvv7WUzpdkRWe210lfDmTpOhAzV0CCG4h51qVvMRbb58Scc5RzMkqOyrtv865PwuZnccdAOU
vsmyelz7V/hufwpvRPDeEyYwZwJE1j7PM9qAsbt6S9gtKbzRzxNEqTH8km0vFIN4zQ5z/DSc8w/p2wLJpTQpN4ajMTuCcCnSNkM5H75i/

tW6WJhOUqnXzMVApjYjU4Un18yjm2/L+LMj/hJ11VPbhGzMp8oS52WwF/dL/sHvgLWQ4mnvUhKiX5vg4lBDDyIPsMOQsgFw0FjSKYsiot
4tlti+aMuX3Tou9pJm8rZESv3VLwYVQFqzSS8BEMVVblFwInqu8dHJ24XUIjDlMOY5s+Qspb7HDLXG4KL4+VNoGePsJAnZet2yp1U9+5z
ZyavOyrbvvU/4OrScEEt5gMYBaKQ/Y7TuJRn3Mem3SJ8JypC01/OyZkUi0mrKobChH7eDmU3GOiRXtGLR9LkJmnExmvwuIrbSDgjynvKo
RULYXqngp2lPUtznC61G+BGIuV6CCiANPboni58eUag7snIRMuvGyAF6W+jF53JH07Bdc0aM7bABA1YapQVUD7iLGv7nIaizp1Aht+F2Td
wF5fmW2LIgCDoDsUtAEhr6M4MmEZPI18GhK4aNRzjKf1lZPktAmoGxchK2OMHRkQCCbT6IyQKL0S2tyHy4H40rsnmR+srr+uCV7lEwH
NMDR4Xk2C/iZa/i3dtqDBN1WInqftAar9I/IYEWZN+IVCd9t4WUSwXOhyCGFH1iKZE4KUFuU8PR5PGSbaULCHo6qGuvHZytGcxJjDkJfSp
NAsykI/UFluax9LVhee/7JwiVc9PlBIpxVwJmqhg5vLXktLTgQFJ9WqH1UkMqNzqKroD+MYtVfeApUjI5vXaLrBhOirymRsx5VahJO3eF/wE
54gpuxOGqzmWoM9LBlOs+i1MZwEwSC15kSGKLXJqvOTW++3fAwGb+Al5Jk1lnxqaaC+M/eeRQFTzpksBWKUamV3SN/v86Dyzk4Fzsn
rgaXDTa3dAhEyg8eCPlakWudzx0AH3wo3lVfIaq6ojfIJfrdMCN79fLRyiYUehdqKPT9nyWggUqjR/FbDbYPK3/5OdWLBlU+G19Vc3N5ntF8j
OU8ydvPhxY+WjUPJc/3wpRTFXMECswfzG/pXhtCdnkPUEn1OHtlr2huVSXbl57uR55NqHQgJSfI4sm6dE2UIe3TPEEk6b+FZysTU3fjs/FN
dOOVteUY+9T48k4T9DR1YbjYpNATp3G5TC41cQ7sfnJJt27khcXQPx6RNGH4hLlbjhbsVDWWLvMRh0MNS+1nXOd+ymPLWmxn2MIhg6
83jmaEZ97/2X5WmNEd37YSeVeF4h4H49woaKkSP3PQ5XNIun9UWxNj62Ckvjtr0AxIm8ldQPh5lnqqRPWRWazxT+IbNBNplhk5T0IoDK0T
pAl4hbpRQpqsD+wAiujIh+fXz8efKOY2pFa9I2KvXQxQP2F4v8/9N+iwKWypKtRrUo8rm1FszlKLBPhmi7dxihw4TCHZTi0/RB8MiLAcgPZy
09KNkuQmTG5HPJbPX8DVGrU+Uvlt4GEhl9jSB8eQklcTm+0zR3feN7mB5Op3jAA3rCM8ZSlawcUbERUhiJudWYz5wfTBSbQu3HrL5yHw
NwQvxkqDgtz1yqTvyyTDtvy0mgn5jn3Z1UD7wcAFKFNza6/hvnm8G9Dmt+garVxeTHnFfCIwroOVtFdB+2rRRuLZyn80Y0/qTHDhF2Q79
wdjUdTS5x3Ls9BglYUrqS7obN1nlaRO9WZ3QzDBEhx3GMvR4v9iM2D1vg1ZjcQpUh1g+UwW5ekF5SpA3I0i/gUn7RKM8BJQPG2DF3VFC
Q9gN4WcL7bmY9ZCD5AYaH9ONkTJSDS5FrNuWKJnsQQgB9//BkCq1fWNLmvqSL+C7qdf5M1c6G6q1Hh9+AD+wPksamcWUnPflMDM
/u3ARoUPWED9PXyEr3Q/6E1uLIqKSF0UGsysafMgFOCbDJfXM+lCeq4q11udonwE7o2Mxw4iZWsN71sjpiusJ0h+Y2P4qNyUlgrclD6Dr/0
MqOD6M3nPYBAeRuPuyyCpWT4gaX0/zd8p5iSKjLqW87dq65tbw7CiYGATZ8OtmNI8bHC1wAS2F0NXnixHOzXZjsz4s77TpXjN4wcUp+S
xqsjZJ1DC0gSlllZ5zuh/hPLAHzHQgy507BuZlm1fyoNM2HdfbBVezlEwTEyyjcd5qZOgxOukO0ZzECl17AMlC0IkQfVWH/S90yZXxxYSt5SK
BXnzYaIwDnobr/hAPS/6UY7LC+4l83/Tfp6WFVhyfa/+vH78NPie4GqM2XLZZ0I0sXPTE49fhdaAOsu739F5j5hqwnCRAMK6Y4qBeEiKSu
W16J59puhPhn1xfio1OLmfwSkGR5X1RLkJprqaBlqjVSdwh510WQcGHI9TguWlPnWv/SaZFCB89dbqQWpS+SdWDfTXQizQM4UCvPF6P
FWFx6ZN4okcLzt4KaSEsClffPID8ypeKn5lHlC26PxbzGT9ubTB5J3pUmEf4ljgoB6WwnLUgXTz1qOeSWVX2PAp0dLIgxrOb4gvDPFnQbgu
6yHclRYwYNe7cfepcQ1Ign2IRx3Va2Vvx8LONf+udrc/4miRM6BtueZKiAJ+XmP2yd686zHodY2GeHcriaR4DcAjYoCpo0gFI6NLTn80QhuZ
1EasIM8IlEcXWZDgWx51iNgR7eWYPTXQvEoT8uivlS19ZXx74EeCg96fyH2qGs4BmH5fJN9sHkWis+iPnHbxOMixd0Fc5xTz5ehyBbzD6Ja
1fumvIXYeNqN1rGlzdipWOKKlAIJd+Jsjs+ZnBh9W+9+QkV+ZLJraLWaE219GCmKnjF5sWiVylDjp7roEBCa4f4uYM4zUzd4Hio5YZuLee+U
KiOzUTj+j0bol809DePG1x5eGf38dVCiZ0MsEaPTG0zpNMD2bsNlWQd6cJP8cTuzHcku4sjmySDN/btYAnz0pqFI4Tzf6Y4rifjuLYwaTUHKY
dwGUKcvdPbBlm/0schCa+dWAEAbg6pkn3ipBkG9VjlWM1z3SAwwYiNHlw/aJBjWDbfKYfGqAucwkTLshSB3mNHPOqyKWd7HbOEnYTX
xPzSD7tgqlAvGow2x1fv1Hz/wNrO5oLIoMbZcxJoAroqb0xZi0W5GZKL6eqWsJI6FAv+CoHyYC8YVLsarTUAB298XJxl/j1CKdqn+xllGyK0k
ruaKkzP2oGBosr9tX3hbPRRoiYnge/Bw8ibuSRunYoPLjgdgwKt/mahNy6IWlrAtangY+oMTWNzt57L0C78eXpXxKjeTGh7Jxm77z/TpnQw
LM9gwQNnRCFzAUio1H0G3FMf+eAQA5OLcz5S+Fa/AeEhLTyTlT0RPUk4Dm+zZqNAqXdbuYQDfBGhJwwVXsHGZ1QylmVOL5ZqL1Vce
godie9avl46l04+/cbGTitCDU0YFMML2QT9Dfrc5LukBk11ncLz6jlmTlIY+bEbTghzGGEp9Z8JTEBfyYnzz5GCGxV21wZGVCINaiu8AqODG
NSAK5Wqs2Cm0Cuispt+4/wm1QcopAtI1u7lK9o7h+ZRIehIT2MZsndZsoQYCZm7ZfjxxqVaDUB9r4V1R0ctlt+n48pKgcsbUITOV9wcGZg4
ejEKNzXgbFIakcqqFfzambXxaj5FYlTZ+LL8sPSJfK+wVJ3Pe4UI5p0/X+1J7Z9zp6hl3/bFIqYfCmmlUff708v4tTbYpsqnrzIkr3Iquen1/+Pzl
wRWCxSG1OkMeDnn3nXE/SGIXJJyVUO37Loa6WGrn/Gq/wiKvTb69Ptb03Av+yXEydyfvO6xUlE2uYTclIMgF/oYs/NqTGknBxdK11b+t+5K
tEwBEe7TNAqUXxi5S38UmL3xTUH4TeYAMz3TNupWZ5rNCh8uH2AUZSNFnAvpbKRsybl2JkO45MHJo0YuQF9P3xbHe6R2slv260oNQ3
5CO4UJquTeEyxqawbSPiY41rgke0A/8AV/FSGaFhPsejaWEiU6MW0SFWiOEszzpgb+XIegLucwFdHAOLFkx7yImhfw5M16RYMbSLhWIL
OcHdNCIUU4+XsTQVFok/eqkrZQnaRNsfCnLCUqnEZZKsJLE/IeppMcEKtmolpdhO3RI7tcciWYX8fplIlWx8g80OogyQJVpyQIjA/I6F5wfAjo
j9X3jY7+Mev5rzBTffW2ZbCgbGVh+wBDVC3AOejBluxZK0uGwBRLLKl1PfdTevdX1YguL94APwxolDElB813s+dAvDODExDBypP7FCohCQ
NVmJdeHVpGCmjwsv3GBBnxG5aNglT0fwMkBF1cGX0tD0FRavneRfGE6rtUmIBkhKSMzOw7w/u2U6QsQjV4+NRcU4IjLIp3ZCJ2ElAq6
WfAj8AS+GAYlllltA6BKWFiR3wKgQiNdsQD16y8tqjVnakQWGBvPQkUahgFs0jPK8BCWEnFT/UuP0OxbVncMF6HDzghdl+mpZkpPJhjU5gJ
1bEbfNlgg7TA6eBTAsd6WkBCoDliZkkeVsHCCgnKh5ddtUZJpL1ezxZDoAS5zoRzuECGk6QkAYMVPTCT5iEyv+67+b3hsJuybxo+uyEXsJf
KFDuAd4xVSF5Un6EqWdF6sTThl7IJsQFPeXGHvHyuPSp7K5FSy39zpt4PPch69wRz1tppjEc9nzPbnOIaHonIkAE5wOn+cKmiXhTC+5j+4
mX4SKWoENUYmFa9wa8NQVzuF2rbU3a8rp1jBB7MZhN4zVwMAKe+aPsICU+LITehbvbbE/y/C7kRnNCLLZ2VxK5mdoQ3lWqIz4rbS041
T1wLzq+uGiluhbymvvS1zjM6pusNsrp8qM1yiVJrl4vBY5P6TSC2mpj/nfzP4md3s8c+EYl1MbWrajrhtn0YTd0gx1+f/4y2yKQroHlGHB092A
8KKikFAQhfn/LOVldT1XevBC0O4QiXT2bl5vqdofdYFn/EhU9ZmES2pE2pW/JXCEfB6goetnrw1SaIYJ6zWld3eGQaJNdT6xkXPENoighHiv
3RyBsRrj4LR0f2lSgu1HLHl00xr7vIMyxWWXm3DWxsrwDRJvPFBVPYwFV8WgAVs+hrn7Xgmsa0M3o+PHcWmMcsp0uwezTTSw2d4ksjj
oxszreFNuoVzE2INC3L5pjp2eV1KrBK3HMF1slTs6u16VY/Dh7syOBqQelOOR/OVfBkhLibW5eWwznxfSuaTHC/juz2GHJBjkVFBsONAVV
dWpD4zLiSX72dJmtkXnzp4dlFj70JJTZkw/+596SUpMv7vuTsuH5Net69kTMT4h8y6T7g5i3nDLNT4OJxscAsVZXshZgxKUuifLXYArOCeQ
jMhXRMP7i80QSQ15CioTeKOcEFFvkOnNGTa/aN6QMEY97VfUvDZw942ZBFtuFWUQh9uxoluoqVLRqayp0PahoONriHdQTtKqMD/PjYza
iGJ2HKlEkjGtfidwWPqJiBBYOg+cX2UOufQcNvfQ9Y1w+is0AVaCyu+2CqFmUKzM2riAf4ux4NYX++n9DwphGmDQgBmI4tIkJWXB8Kx+M
j+dIeehHbUPXh3a67Y2ZcuDTpS9cGn3x0OwBIQYvYa2LMTts2jIVZSYcT51x48Ubet6GHohPXLiaH9nPupXMsUnet+tN3C+t68iYs3V9mac
kCsM9OI4GK/v0uZDJlc7iQDNst7wUV0cdRwlF6kGgyfUtHI++eTeTPh+WMB/gIPq27G6mE1RoODc/xme3+gG/t4AxcjkhSZ+LQ8CD6LVAu
IjKKQm7P78L/X1yCjhALHGC+6Zl0TlvboXk8yKrrsus4Hus+ZNqMDouHYbmEJN6sXBQjCNytYmiwtqh61ddhL47CeJkideaLEQHPwJBE7N
Hg2iDQvNYXX3pmkyAXgZ5BgeT5XBWMkgZWNj0rJBDiWjS8O1lHR8LL/zbf7OdYYO9rPTphvefYVY+gRXHADbK4uAlITzacYX+EEYt0I/oH
07IvkHMFVkE0rh/GgXhBBzN47weC0Dh+kei3i09tBUWbCn4AHMNEIuE4Lgas3zyN/939YGEh+lxU9wy5bSuqr8FzhD8DQkBCNlb1sPeRoZ
Gcub6pmeT6IpWOM8ArRBlLpZT7khlQx+0dPCTaXh+FLQ3vIO/5WERGureIjptyLrVp/yODZgbHpjSzoRxBBbho6pmL7pbQPWCQI4B9sZ7Y
17G8M7TVh2dYJf/qOjHMGPfEUiM9KimtXF++YpCpZ2p1pLbWAjvGzObjD6tBL7c4plQYlXO0eqepijUeOry7b/Tg/titKgTecXZLhP//UUbyyHix

qmgx4tOFO8bAEGFvqrbFadubLBZv2GwMlD+42+pRqcmllBvpcln2npoHeyNzuc/2mNOE6sxKuazqzuPqeoRAV/LglfD/2dJfWWeb239s0
WU6B8oirdXprBmdi2d1rRrgXHFn8fABoJhSWnEuXAuBCu0dokf/bZsSDp93i9lJtKI/W0UwwkGh3pt78dhCuSHRIavscThCMggYPTChkGL
eZ0Uh8WbuXPFqNfvDMMoJGlhrKMZz1GT2ExNq+buKpSkSMk4ODmdGE+SihM+BrBTHAJBEaM2yUC/KWpd4H3U2Oy0fWMyh/rTudQs
hBH2gYRGYd4Hzjc594yP+owMmpGSRmsqqWx786VudMqa886KKfxYfRuJ1cFqxbL9yYgNRmdZIjFRETfH4tdF/EdPRYHxkSVn1O0Dt57
62SMZ9U/Xs8KDwoRcuDL1m8dUB1qjFEmu9n36ZYZUU172y6okfcXmDsXYWm01tCYlZxfyVnuPdIO6FJsOEcZuMoTV91sjlEBjP1YlCqKp
cK9DyxZ6Jb500yCqWcl35j1Mb/j4C517wzLNyDpRtX8nWciKgETIuAPsE7hyK+DVb+Y9S3S5LJehsDDcz5bCLJRgGnRrEOEpLv4VNWavq
qDXP/01Oy0nfRc9KWFYa3qTQHGbH2czuR5k0Ye/D4X/osOZEh7HqTFIW5+gLoCGmWJAB3lbLz1Q8RCzvEJCRyjv+VDW7MDJu4fsI2G3
//dsHv+wn8Oz/2NCL2CTB9Gi+L7YIGf7TLuXRtNr2+eTxUgZQrWzP8Pd1je7sq5jP+rUzq3ChJoZmLnNFTLHeE3VO2FytsTbjaK1akIEzCH
DZWZ6ZJi8sWP08Yb8da249/w6nhIwTLlGK6PHCPZYrZhNOcHcCF3X1lb8ggKWrjtEXV8Ykx2gIuuJsl4PhvkKOfCpPW/Ix/MiaaF6U0jjjiDU
s+mCktdbNzz31g3jXpcTwOO0b/Q0srZA7+q7NryW2LUYMxyUBCUBkKZ+/hwZIF55dR80xSf4+mNj9S53ULVgzCFGqWXOaHUxnWgq7a/
RRa08TVEXzZvAmWKLFwj5F2lgxzlxuJAegbZ6uCU/h9EzN9HtJr58ABVbnwH1lbsfausJXV4Gipe7bQDBf+md9qmy5wwmZ1ICtcOy8Osji
asumiAGfPbYd2XANGY571CTR6l5T4dmrQjuZoL6f5X87Ru37P3O/yz68rIJ4p7FcuztkyxCZPHegUgx0C/N/h+bldL//lIiWJYaPzi8ngYn/0g
Jii2dBX1VBYTN5rkk2eTVIcItLOrrnNrRAo5+SitCcXIoLMGK8kFmoIRhurC0YmVTfkxbCV6BsbVN8P7P7VzoslvzztQZK0A/aX3fiJHH6cOt8
3oJ6puRyHGrOo7qNhQFhwYtjJfAYGJ/jzY+gUkgg2TK1UQFch1syBZe11ogfdIeNhaFMuU9Jhr5Qik2fXSHdzmTww042ICslZYUUveHhz+t
eNGbIHdAvXZsrBnhbiOHeE8S/UBOvJKooL61DjK5aDHIkGsMVw6eQHmwik6bgIb+T4IJjCjBSujG6zYdYk5OJ0wFTtaC7yTv44CD9mSOlD
8sd7v7OMFgJY7qA0B+dkQGIOkrYIL96K+8AokqNswhB/MQ3sbo/QTcL+4McmNcFqfQTcVp10V8uqXGew90UFAWTWmOsSVQcORufpv
qwCrwG4li4SFpDK7M5DNKRxKe06ZSIHSfT5e4Yyokl5wEs+yiC2YmJQn8gnuw+Sy8nRLG3naUFTS7/iDC/AKcQCqjV23qPB/SmhfauSnSh
r9JsXsSChizlSo/d0vjcfGmCkirUP3ewtzDEjppGMyR+TFIm18WgendujPIgmGiub8PuRfVDWRp1WXby1oqEysP/UxlyvoahPMk9A3xvm/0H
21VmbaldytrZ3cGBdT7e5H34yxjxghiPZcmEl65njynJcoOZ2KNJsSKWzPazj1XYCaRjDj2x8dmHt2FefpfM47dOPRFBJ8nIgWsioFjfo1R0DA
wsPTvV8EPDKGplXs3cHwCnROHISdIeO4V25oaTUk4ii6ZDgg+3tHEZ+RwG2Fq7niZO8wlJ20oNV515qEU4Dtuur62T5pdArmM4SpsDTB6
8Z0U56GxFbBQBy+7rc7otedEARKYpdMiPs/J5ptHRztzmJjQ3uoSUReUPPhZ9ZeBaHoI6iSR51xbuYWq9m/iHHFIHMRBuQRS+CeCb96DU
gEAi6RN9ax6zemx+JckWoACWTwOrJz1yPQj5eof/8nU8ymVH5LoksfWVj5usVcez/qj5zTDPqltRcQ8XhvkpqQVMxo27hPiyOtyAWtxWYM
r7nMOJ8C4b2MQCYivYx3/bw3Wo9c2BfgkW4YWifJvuKavyd4HXLtqOqJpLJKSzK+KPzoulaumgTe1wX0B5jc2KlA5GPDbafwrQTx3REJ
DTcxL7ZZl+qe0ncqKyiqKOglxwPl0QCaG0xc0KOnA7MvUJLNI14+nXH17tCJ7Xlq+ipinzzxdG7hPJJ1XH3Ef6fqDW0B3MLKoIGpoi6Cnl+m
xS5YG7A/BGspdul8Pc6yvFRm/ZWYy+omaxRr1Ba5CeacNKIePiZtI70MpGSul8FVidSONi34cGsOUAGhnyZiaqZGs7OgyDYEOOutum1BiX
8w9zCCECwegYbi2oOFTmcnHVog/e7mmGE4MqOIHKnXptlKxaoaptUKx7tQKhjNh1QbZsVEJKUTzpmwW8901yoadqJucI01PASiaOjn1s
VFM/lV0V3uQ2PF9ghyEjEnaAD7TT7zR+NV56npQXwczZIDCPZYRPwQnaCFmCShg7a9ApkzNeTweoiJ1LVJdhppJaR6K4kPUpv1HUu0E
JFHk0w76ZztoFFm1oZfWv5tJr48KwR3pzLbt9SAuBE5FvxUy3sz4EkeNfssFV5kYkU/uNHyznWT5o6NGRu9MB20a9REeK1LTAnY0JTIO
XU4PSruw8285RKwWbZTcP8ChBNqNQP5un8DxfqDVygUsDAt3ec8Cla3kjxQmV1j/Tyv8ayTTc3U2R5MD68BItaGMkaAuEee1ziZIdeUoB
UcDs+aPSYrK/VcvDoXZ6hQQ/AsaC2vd4QGJOgYtUts1TMoBmSk0rXqzTRokVbovvlKKHwjBtqDXwV2fz7A5cLTIBopffHV2Ia9kyUeEwiVV
5Sn0hLiZXO7anX+XfUsksoA1NxLrrAxWKyNOeAmLgsYvQe0Jahr224O0FujD49NyWmIw5xfBMAIeTHSUXlnTR6/rVC6AYfrM05z1DbpOT
n/KfnxR8zVWUyTDS4QGxcLTgspeNOlNQLDpK400MZVTcVMPL/eg+DEt1LEwiog+1LTuGFI/RkXa8ttqECZ8zozgZLUGMiRoOK3gsDpKR
8313zQ6LggkYaBoQ0m+jsswUxQ4ENFeGwf1uYqonGSmblQgI11Y6GzLhYDXRmseViI6oBmrICyVFhyqF3nzl6Rk2dyB1p8jklWSJwckc5r/
crU3eFmdzXS6us6/Ajbm0aP/0+d4bkVcYiCzjQb/QXdIk68IyN3KUyV3mxf82DqXtcsj/d5WuKdPNV7AykEMMIoP2byVagznkqilueKKzVQV
SjjEdS7G7d8oFC2NkBk3nxNQIUkx5rJyoIDb5HH6djIF8rHBAPyoMKtSEPX7iMkHVaKCyDQG5fIN5ADRLbdItFgDXqBYBUOvPt2+YsuEr+7
ahFQFCGIIbQ+z1lZos7pShT33zNxI1Yb9MO2afD/H7CImfH9ufYxmI2urPEMIYKgoonSdDSmMkdutzdTx/CQkYgy7gyWLXQEE3d402mOv
/WBGfn+yhs9y8MM9hgY28il26/eMvw99QEJJ7UmSx0tjjJUSytrEK+D3GWf+QFfGX+zda9KMVIXnhxSf0mtEcUjfco6FySC5j2bSosufJwKR
9Zm+sdMQ1b8YZUTfXpJ1btlnOMWh8CVq2nvslikcX2U+ibXrJxLAzebQy0Ovh0SpOJL0hmSQudcemQRCKmzry6kXcPADdDieG8j3pF5FK
Uv/dt/1iUlE0z+c4gSjYxK/6EbkNzX94E/wXFFq6sNTmjolFO7nN7DjlEaWmi0IJrw1wRSMc6/t3tQaabRG6SfrE6CBYrXCbdzMrj3//y6UAfgO
QTuW3kdb7ETuM0UB+2TbI2ElAbQJSL6xvnzeJsu1VaOvZmu32z7IocrR8YFQfD6+yRvMb0V/BIu0UmEX5X4elhRVdC0lwl1JjnzImnGh/s5
5llY+HhuT2kGXM12rYccs8dA+ZTuMR1Ol/3T2mzgyXLS4Dpj3Zy5DfZ3DjUq1m+vuSZEXZXNorBF+BJCOv4TZV9C1cGiS3SqVjZd7ZXdNM
0kFEmdKwJI5263olBqz1wAW8YqNfCMP51zimyDq/ufyfIDJJOp4EuqfACSduUQ14p/dR/l6xeZizlGcz4ks6Rgusg4l1s1PcZn3FWAvuIDDSl
cync82ukbHIfmexawvE+Q76cQWmcmqSOh+mWFj+ugISZZUP+1GP9NYgbFyEaDSjtGSU4mtm7FetK+YwNcpRpY7k+Jf/QGocBM12VuH
aIs0dBovcK3QUY9K+jbF00odRUt4i850sOD5qDToJk5rrEnchBD4lQ34EOrJG8JgEMIo3LSdoyNsZjwco/ruY7PYDw+RY4mWR0xc0LBg/Vo
tgyk3291Ic3Z97q1pk955gtSYWMoyXA1n3Ch5cSRcjFd+fOwya5+Cq+jMjRlXFsl5JIabWZ+q2PsL8RrMNeeyeSWOBwBypCFH5CiBERfl3s6
4Zpz53GAIuvxyEl6W4VYKjT93gzy0yjjVyjsHVDW6JO+0JxnoMGoUR29E3n2zQr0G9ytMoBLxoYohq2FTHyTaDNESU7Bjhn6Yg/6DiwQ4+q
xuz0zLiSigIF3FgoicL9I0xH97m24Sc3d6xgJXpmEG3LU2qLkAkUbLqxh6oupr5JhExd8dMMK46goc0LD6H6m78Jx7QfbuJtzcY/3UIBdtJ+
eOo5k3klnF06nNnvuwEFOOailh/xA9b8YzzbLClHNPJd6KS0faxDt1XyxldAUOEhHRhqAUlySqfsKhrQKJ3hi6gJmwtXT8jEOUZ9UlKoIm4lS
/ER+lXoy3wCeSN0mqx5ZmR/YNHoAEq2LN4+h3zLo7pVUY2QxRls8/m29nTkigv2kCl8z/P9cH2TNMsASw5kJklZA6iFypKtQQ/VsvuSZO
GZjs9zcjb28h93IUo0j21nEtmGi68CRQglVlHaXFGmEVp1qxd9pUaVBRGLX0ODrzgCX7smC9JUA5ZOqViDsghFWMpQ7Bps1tPncFOZJko
SGa+sUzcPG0ASIBHHbzEE+HpgDZjvpVy+aZbHIhHh49/GpJtkdKokHs/5udFSRwEC7hz6s0+I5CLww4EpkB7sHGDjSk98Jp3m9Vum5i/J
LVHczuIJakWITtV027uLba0otVD5Lw/KqyIzYT7/7DUFvzMzE9RpRiPTSv403eriFDq5hzq85EGxyLd+78ZptmTRM8jMgEOn5aYuer8GmbZ
6jbMwm77Aqem9IDwHXOG8iZ9razje3HObchnfDlRsqPPKOIUg1PuR52SpyeeOoGGPVYZYBLfS0g50VKfKhhfJ65p0jpmhVxQ8ch/U8KVp
rDtWMJBwOTwbWa35dVhy063e3UhXhIuzzPTaS/l9wWGv4Th6FfbOcOdziZSEU+5FCTcXxY7t+OYctV4JGndPllEogwc5QesQfsb4hN/Ug
G4H2DkrBgaWrQmuy4ejwChM/m0kdfVBOtpuhQnmnvrxE00Zv72YAyhGe5NYJA3b58C5LorloZpiQoHgK/A43GS915sof0Y4160OiJVuP1
diMmXXAdKLij+bFcBVnxL7FXsgvuCqkPaUfMpBWCra7c8PuqprHctijnNR+tC8Qvzg6wQdUZYFEPEQgzr0cXBIjvKq40L6Ejxh5w63OLs+m
chrvdkAFohtZ9MZv6k1twkufqFzg2BFnI8AE/DML2Zyu6hiba29Zw75Yg5KDl4d0YBZrfO8hy6TLWgernSkUILy9PeaCdFwAqpm0Qu9fW0
MStnBat9xikia1rmudJC6e3BAT+wWNJvop/zV8P7O4Cow24/9wiTTqYBeq5JdUyqA8GjtOvHCaIiapw7DaOpR8dWlr7oMVl60yP3Rn5GR8

x21p4W1J9OJf7MwaiqyJtzr+jjF4lp7625HE1G2aO+tTS1Kpz2t2mtiDaxlX1Y5+9J+gHjvdfvFDQpztdbYtYd30QZSjnAFmi3KlolEKz4ukucvl
GdMsR6dU45UJQ4tA/rEDgbjHeSCdSTVFzhZbzxULySRrxBJLpJYrafQ1hCTM2ypPhREMq8XkoQiSGjM3THcdhPRi76uRSFfrdMUPvIvSr9
Y4il1aEFmZhl2a/C0tiOBWBm4AqIFCtLZSS4A6q5u/a5u5SoKis/bfih7ZYGBZCiyEB4lEh96rHKrtbonwY2hV9IF3VQOMrVk9xQ9edUOE+Rf
Q4+AY8dFXIkYb5piWdE/9gRSH9sS6tuPMS7zlpDvWwachWJ1z5tU2kb6LJ8luG7JEFbXbn5mq3XOdSVf52nEzlMXu7n86zUngafDYwuIW
htFkDAM01YpewiEaGTCx3NVg3/68zr6XyODf3OVBYkechOqQyBWlnjXnYFLGmyihbgomz8HMagPutDHFo9FzGTAVDFGwlnnWNzrtV35q
LoOhxZzrw9jaOz2yJosTlk4HqqjbWxLwUP3aIkzTpJtiPRNqOGStAZ4+hm+lnUvW9knCE9pJM7OEmYdWVoJyIP9NKiRJJgjYEBe0KQNq7
dKl2mJl8wjwJS6x48U/2xFGiP3nlhd5ux3F+KAHac3yT3bQ4uPq8/E+WvNiol/H6U/l2aih2pOqVxmXHzlPP3izaTwde9wciavw6UhwcMDAL
6y5482nv+EfEjH546sevLzmDysUBwbfskBiAoLqrBfgdVGiWdA+Rvojz5pxedGlMOJAZd4WJ3x9nlVYtcVXSWCY033WAcpZr5Vu4JJrEzU1
nHXJ+jT7twckU4sruU8KBlzVk/7dorvXLL0f8/MlXCuiXNlTHkXQI7bG0iURMQSSjRiN2/ETDXqvo5g4l/ZjMr0QBUNwW0HLmqHcLFQj5xUR
qOMectpPJ1B1HysGOQGoft3++MYhl6wAw+nzxu4BSWB/w8TPdDyMl9enOzM5+alh6GVITNieaTraPPAWstw50IZBmRPKVQPDFp7jP3
msEQjjETootVCm4BXJQ6TcigQUFw3gtBNqjyRZsCgosGwyv5nLJSnD4hejOf4XOCJzgTSXukwV5zpPMgqT+KiJil0rNlbcjlHe2sfhQMq0EN
dZmWQrRIYPNg9gnj6py4+Bpelz5/YY3i9oQBLHxJHVlaNx3TDF2d+qXF1kQrp1FhembYwtUvDEGf1tSVaNZS95Xif1gIdUKFjcyehkPlbHj9
WYxuwUIO1htmj3EHwoP7bM3zfrYbCufNVSldw+rG79utJrpAdytwuUS8y66nyiSzeeVj68pXTLDUcEb5bofQfffFClUKw7MgWhYg7Q0U4M
gDVa4wWok4W9Ofi05hpk6mxgQlK1RW3ah7xGyOxhl8L++2isIkcPzRFddXd07BkueZbwF3qwnh1nY9xaaoiR4u4ihyLs4P44gLiTavcrOsyZ
g7xhvjF1BHrtehx1JtQ7q3jT5cRRFiFZyujFVfdnrrFKADpCBNzWHjv7bX03av+8XNqw8tCwJ/9Mx1uA1vAZuf52zt0qh/1Kgh2WcbY7Zqs5F
ocsI8D2JA4KopqKtq0WF8bkY/kTxRDpiHVHGxDuDKS3flZ/u/SZI8S0AIMPBuM5ZLLUqvfAEzsLfT+26ZMFroWYUWNW/FrSpuy/V2zPOH2
RSgkRKcmzgM0ts5s0bGc34gLfOJq3xa1/gDhQxh1wCDF9h/6CFsvY/hu3nn61wZZGt6cEJvJ3o+vj1hdXFUG6FSG2TSEbtFazHfcmm34o
2dqPkMJ1+yTFanDo7xtXsau7+KmQ1mxL/xRaTiH3CYRlYJSJQwa8nYwjJT2uMHh12Zc4wCbypA5LDYs3nt7KqTMiM66PeUtplaPWia6e
cpljyojQ/6+7N2rIwEpNYKFBa0qodVRWuBuTprIopaMWWOJEtp4KIEmapDPZRzcln/lcYJslo6OwEJs0CcvahqD79Gz4BRtK0blcBrZAbK10
+qGfbpWxhVel0rPQHHMOEKPGbBgs7UI3X/uBqdv0QNt2UtONn2M2vjlOgXyaowPDtF71UwZhkXGIpYV9pZpqd4Ywv8N/AGfA6QYd5OP
aEdb84JVk6IcBvKk1P037T6FbBDH5rSUTKjB6kANTp6tA+Emp/vyou8lUJSk5n0jFmmK2fOwb61VqoJUA25rTphRgwyew+EWs2L/umaO
e5aLG4ltQJECTbFolwV6QYRzQwjBJOKo102uE+zpRhS6dbUzdwR8y5J76IF/LbCWw708XTZ3GC9CS2bOTtn9zA2CWS5xaPSf05X5GhU
1LbAyiFfl++iWafq9/i1umctLFr9+nhslsjpMVF0T79n93q12F1Ayo2TZYOLoSL5bn/Y0We+acOUUOUquwAx2R0BYjY+qzMuUQ1A0TJdlhQ/
xIF2+G60kTacfy/5Ss48xRRt8xQD4/kSufclApe2Aqx7DJc/IXowdxNqMJ4mDPTwL92n1aPouz+cAXnftZLmq3MIzMi5y7G+j/m+xeMhHHc
M2kJ8hmyuM1hy91G9zCA7XkWCCiV5Tn1puH8k9eox90EmjBrWSBKYY6eYQ98LLIoiiySs7WgfTuxoHz2vFpjyi05CAPm4eVtbj41CveXhzh
QqlyKLrh972SLG964yqDsnYb1zbbQUJK4hPDrgO6iEa9g0+9+aWlvKNz2weKtK7yirMQJRB7CpuMZOJX97MSEXA2H1uEA1ibrSd/wvgQb
tvVLAw+43TZn20WO8ecu+gj7PzLtNOlFCKFaRVNoCJf3SRWSmH/yrCMjsKKzEjK5/jmGcGWNhZCW4gcpazu6sdbHt3XjQ7GXKh6qpt3o
D1yA6bKK5s4kAHRz8n47DeCl6SQC7aQaGGYFtMWG4fDEM04w6onD+IdPTs9DuDvc5YYSBdFi+MlzgxIstLedXBAHS/IEupBk0WqGD9gjE
3mMTWsok7WCTZlUg1IKCpNWEPEQ+fROjsy444cl+7xkK7oEM6dkGVFtYTuLA68OQJia5832GRDNv/6ejHq4B6Newe/cdYGQw5O4ZoBd
2vjl0VMr5N2vw+bx1JMF7HnmCCTuTxzeTsrrQvj5yF5IQcVOTKSp9CJlcWdf8HXWYPBpGGVO8gwZZiMddDVOXV5p9o97rOSuY0hWiNE7jL
6gRCkRSeeQVjlptz6elqq/VAm7v9bUH/Pt87lTNHi++FEiLVnp1mKdyeCnjcfwW1Es5ATGjlPOorZNnO/IRtFpM+3QEopqEMwVqRm1znvhxr
nskf4O+93WuAkfUNDym5Zx7Un6ayqo94m563ZGxFAmtDY5Doj2P0+GGXxgF+BzVelFb2CyC0wt/9end2Zob6Q7NwC7yEuzqvCrJB7P/o
JOvcR3MCzjAAjPyEtB83YZ4vkRSHuqd31U0odbX30snyMun9a97wruaUaOisBdga41YOiRxoZ0VM5fbNAYm9Du7kY6AAD4YzheDVkeym
gv4qqKkrYJPFWJpDKnhxtq161hiyimL7Xjs9wpyBVyYhJnn9gQahn2RCwnDuniiEjnkIPTDztGJYYcBS/QQ4tgavZJy52Q6dsIeBE7mIcabxb
Voc190CuEA3d++ZJPOYK7wOy4HySllZx4cDD9G9CtSzgQePzt8a/r+t6T7ZVYf2XNXMyKgpRuxJ5PgX9TlKYmrAO5JAxaCNg6G+3xJMP
yiJY/7Suw/yM6vQGrqfFGsyv5aikv4upQRDqRoQuyki+1UND1ZDlZfamPeZWGYLCehMiFX1Ji8tqJvt669beP3N8/k+IY+fBO2L26JupUbg0
qigy2jvEQOdq/rS7YleErtjTu6Ed2WPAFAyJY7J7iJb3Oi7cYSdDiArXWWCrsV3MPg2ANzA1g6DUwmVLQV7zPa1up67U4N3msQve3eE2rW
iomxrOsadKjOYkC9qPePixOBl/xrIXNgYllRKjEGxRrbkuVRqBQG7LyTxlc6m8dwOxMuDBQpFxlQMj4u0eyCQxI0DSik2MkGQ0/5MGsp87G
NPh31Hl6AWV6kNV0yFjnzvNA6qFZiErvYpmaPt3fUWR328ld18d8XtH5uNA7piEgAGGCJdcAWz6oKrcSgyyC9NU1dKAhFTGxylueOrjTNi
M6Unqc0/OUf0RKs3GhzS3QIN2pB0GWk62xRWrZejAhBh7+M83KTiK0KVm4QsZHxJLU+K8/72E76SCrhGLIyFoDn/Bnr7C3GqJUX+vHm
osihgYCkV3x1YpHPavk0NK8cgWnlF39uqRmBaqf3HwfOKf0hndBvtGLqrBsRedrJVO3NTvN2fM9tFfKko4XFMOAtkOxEaB7qSyjZ3GDjjFu
OWRVr2ZGkzRmsRf8NxUeYXckDnoereATnE7G1EOPOoN8X2SQLn1oSOv//AMUhNswqpMHkzOr9ITKKbH0a9lg137K+O8ejjRdHPEg2PE
z1xo2JnJ3XTo6HlunduysffFrVkz8DQuKMQFlc2xWZgTbRTdoPPAhBWgBXCU2taZfbRMpMpKxJCnTbJDBS3M7QeTaoch7Y75ZY7d6DI
Ws7+2xUWZy68vONMCPJJtUQzHF6FeZSQrK8c6vvBCMxhPsfwPGP+ExHaXCbB5rCCiv/lWikrUkJksPzM9EePhASxJuJGJAHl/zy9pyDA
QFmNMuQjDENWndHx+prkFy55mi0px7doU3x9c92ZeqgCv4dL9lmDlwmdpD1LyJ2VjRHAHtWi15RsjA8JGx0ZPG96FlTmoaF5IuhSrcMw
dqaoE5Q6Zhz/hRgUVh+D+2J3noYwsiRNhAZ+ts2/eGDfb/lZ4PTlBzVUrPqPi6OH+0GODFSQFFSADmjvf3lx68NNqlnf18L2BceonQmkuRu
hbjkpWiqCX+tpP4ZaHJsgvgEwPWJeLRVT1CA+PrWNsJG0A0PDfpHOT68uyzqj8qQ1Y8w3g2Mr6IJC9FT7AJ7rfeua2V/8JllqAIU/4RGiiq+
4ncQRKh0zEXxV2Zeyr0X4DOQTkTVeVA62ROlLwoOobJNsB2twhaR08/wzNAMixBI7UypLVQ6R9YMVi0kCvDOGpMEwCQwkSGErI6xzw
Lh9T7719N/vLqrkrEdb9tQVoa+J+B6C0U/tg1XSMTj00CsvgGoXRkcluzv/hn9HhBs3UPJVJ7Po8U+wO4CVADwkQpNvX0CM0TQ2a7me
SVniKObMW3i3URmSyyMSUgne9ekp5jpleFuyIIkFgDVzVE2Nkl6ePZvazz3mJgEpAu3p16eCV7MQ0RVg7n7Go9IB+SCVF1cklPnpxDeaQh
DCFJni0RwCegs8FA2qnB7XtJHnIpA39QEEoXOel5lrWgtNZrxz/ripZ4UwGCOtn9HtU22vK5jytWXbX1+niiLTlDnjPaxbMEd/ZpLWZKF+xJO
Gj5qrB1jMXRTWekX6EeTHow3yXNoONVxqIKKcbZSWrqkBPC/YjJUQMXoqwn4NObbkBvioazNe7upzaFPbbZ58+t29SUWLB/DFHh1Hhju
sF25APz8oU4g4iEgtjXMhY6hNTHAeGMEH8+y6Gs51xhraXxedwlXAHB5V0Thjl+7nTW+goCFYO0MmieS8TFAqx8XSkTJzfu56+jeOwblrZi
VJsb9w5zt8ddoXwjcutoEYs+0Emiz1xDr6Gq9HeZPxkfamdi3Ye0XwYq84/y7q1YFMRA/JVdIyLMiJEftRHJ/qDX8PSwOOVzrgNHmRl+vEi
nKaIbkS0UKLgNRu2T7lPVAtmScgmVUxAsk/iltGl3IVvNYXAeNW6X/DqwOBkMTZgGRh/CRyD3r4qaLUFe9d2dWoeRT2mcfU8GP1nEwse
qQG1yDgrKk3cld1oYVL5Jeltj4DoBVhH7moHMqWEX6IehNBKO9tiXUGc9sDMP/Ap2tAQIiy8+AmkKniF09gNTPfr6HxWiLWoAtuormKRFsl
Y1x/n3d4D+ahA4z2Rb8sc36t8pHODBp0qbWOo617VpiC2qIkYxSS+P6BBzPCty/xCMj58/LnC1afz/Y0MgiSeIymK+EozqLv3NGUe2I/dVA
7SXb6Qqzh6SDccmjw2HIT09qFjy46LTvubrhbz9IeY3Lr1UMEQnH3Y+5UFFPnwKkMDV1B/VCagwnsJhh6OqDKtGuDTL4N9vTTH9MFNe

etLb9jzN/wa0SK8fLiWlpbOIyMkVpcNKRecmRD4kkL7bpP08H5WonDuU3EvBLCQyynDWWK4t1nmzk2K7iH2VAGKZ/vz1KFrShK4RiSfFm
4M6y6bAMJzDFDivdzirQmMPO5J1lgglHFBvmaduN0u6caA+VkzOmkOeQjINZHY+TBrGmrcX5GOJOVas4iqZaPAjDzIpiANZH29i4EZZzYA
h/Uw4MuY5iI0KaCoQGx2Houn34VrHUetV3aMiWE8X5q4dp1IgU2qlcCUYbWQ/ybc3yt4ql6XyQmSZclPyJ4kgGCVy+4EQXnRJkTMtb4lIO
0NSZ0QIiLUinKNo8v6G698w0XnwFKyC6Uh0HqOCucxo96U9ZizA8Mzfr/TQ75gCZEOI68XXqGJ7VM3/XLl5txy7B4Z1WNWuWNqWALsG
LEOxc8rGmFGsYCzvKzmBEHQ8uWlLyDlJn1N9rrLUx+yEMQofLhyiOE+Gbu4aBt66dlwix+U0M9gQp6HME1TCEQr3AddF3KXsyitAaEivJK
oCqMd780ZHTMOe4Yl7gqLgOcvvCvuL2IXFzKN0P1tOf7zbhHmIV/6xgI3QYbRd92hYnvwe333wFm/98D1jXNi9/Rc7fFvubITZU55B+S8W
1ecCIXQGwInYrWhLi2xSawxQGi87igvhZ67W/ziJ6V1pYkPiwIZ38YpVlZcOAPCqCwfFopA2Sx8JcKeJNjWr7vw/OtgOF+agsfiycDFYOQSjZ
mzL1fi/rYZBcOveQPaOKNiLRWmWckdyyo7Flph3WvXYwVU8MjdqY2VTUG5QHOyjkl4x5vMFTxYOXfoqZ7FReMQkXyh+adYvIwS8AB3sJ
gnB6hS6VBjM9OTxs+ve+eg5/R67YeqCpV8EGAek75L2sgNRUmG9Hk8C65HZf4ZMjfMyiNyzSW/s+053V3VKlqVoDYwwmQrxz+Fz9LpXt
oOGjigRYTEAISu51jgD8oSadi/fcsnw08ZazazOyZ4DcxFgPELAF8KHc64S7XSlEVixVsv4HdwEWdMsBMOlPNkknZbN/lD87aNPgxZtpxNJ
3hMLG+SFhOebmC++tG9rfacHjG7Euint/3qHmJyCWkqNl20MmG1ZW1X3K41FMPvowjXPBHwx4ehVisJ8ZLQ6AoB7XGE5k3NXjYndeO
5SFrCJElBajXgFcufTyCSALhpTUCSTnjawp622Qc9QyWjnjtlhdmjHor0uA0DxOZJnUt77Bg2nScNCzOUuk+60VVuJzErQrBbyT8DYaiX+36d
gjUpY1bBIIz0qO8jEKoOc2sJDwdKt+0GsSK/TnSW8l3G6bvph9y+Sm6UPI/dxZXPmoveCFbNJlVbP2QvjQI7gQYN8e2NyRlclAxhLR9MOud
bwg6jBgkSN/7BGHKKr8itO9J117H1Nykhqj3bV/zoazrE401W1n1N6BdB59s12qBgHGD4J5m//Ops4ZN1N2Ln0K6kbr+T2P7ZYhBJrByb
KAQVkEV3z1yTmGsjZo/tdHzQxI9QMPfr8Cs+7zbRgb3Niq5A+uBQ0GpCVPn0TnQPpdFF/jf7AaF/pRM1kiG7BmeuG0/wOm+jtzi+sKP+Fu
utAPUwCLc3bJty7rK7zfguM8WkhevatOqbJYqwJiM3kwKZwEYIH+I86NnSyeR7cK7eY2mqV7M401P5MzmPHYWS5uqIw5JVbAwRCkU
BT9Wj9wnIdhExxbiEVbHIZaA/q5NmSm1IK793DKmOzlYpH7ry6kfhkub0A5e+z61433bYlZTHx26FkiZ8/hs6+FOXTZBQE3GPhKU7hv4us/
2clhk6ch3NIGQVo9NSHgGETNML+LokX3BsiWk4TgT7aK2PfVnKs5FbEvLFht0Tui02q76S54GutlQVkO6Z2kP6y1TndQAedGbdmyxSFQw
qkdPac4LDNzwkm9XjEcYvt/RlkLAIYutHqzgZ1yGJZzEq97/exO2IoIPV5o00y/GWfXjgtp3Wj+6r2VWrQKXbZLDmvySqrDF0tV7q+KqC7Xt
OR20Clf1blVIqRbP9PJ5i+IVj3cXcqIsmzU7mYPGlc/0DHzmxTEGaJpW9RCdTA1lm9X8EXOt6MKb8vqV5SG2zL6A2a5M52xUJiHAsoDkTv
+sC7onS2+tRj9vsrm94qfKcWuJkkI9cviuN+ga/Bv7EUGQmaDGkDHiF+qhMkjqWC/YQPHPkowcp0PjXfBrrS1MfBwuJz7Su6AmVDkqrVVH
rTdXjB0pqfYdf+ELWrtgAvEzDWlaa267ujnv/hcV0GC0Nwy++W1TFARwGksngwAcqsFQ8veuCH05a21ijBlKKp7qT5zTfoTCSvGEP0W9Nnt
Dc9ZQIE5nh10BzDVfokgMRYn3DM6SiHlOU8dfsGov2MGozRdu0cdRRPrSc6sErlWXjQakUu0UfCVGJBnC2Sc8cgvWQe8twBdj/CMMZOK
3T1PkaXHTaBzxSvsQGfp/6+DpigBV+gSwmSV7Wnq4u2mwle9lgs9Lbe/MbHUjw2RB6fpSniTvZGXS5Mh9rPCz8R+OxZ6/apygzvSD5lMw
AO7Mmpc4PcpThpvchEylOlv6x6ywbTMVZxoiyeUoaINy9ecKtMHhHpWxsOTrp3QY2hZlBtN9CFE3rONzEccxtaBXpGOvHlFyrqszGpXmpcj
+Fw8zavIzq94GvQ1P/q2gxUMQgifHnweO8gVLzatPP63Pz87zm99Qt29b+vwxeZutQrpUrIimnf6Br9hCDtQP8I0BUvDOZBvVC8W98khCl
XPEXfNgqjVsPK9XmO6T75Cdkxf13eOrACcNRZxhe0yjxjZBOItB4u/WhZ4OE90xriQjopCyCrPjXvn0WLMgae4XOCCheR5mFgHUJXS3dhC
QsnYX0NX52Jq9UpKVFwap5N0y70Or0Tmdp4qIMjWdC0bWy7NUMzEDvEWX5M99XKuY5djEQdte5y7l+5QlHjH5yjbZkdAxRXW0Y14vnld
WkyvKx6jsTkViZ1WZ4ImMed82TQj4vttOUk3sCBNhmFwfpAnnStld8NFXF84RfEVA+FNOlm7kqsSgOn9eOZV0/xv5b0OuuylS7zXOtQd0w
K6j/DgpZPUKQhXrtfYNktPZF1PmsZGBAma8Ct29Tb92UJPdMPDMTWXwRXH4HBdjNMO1O1waXhz8dYFkiJBm7XQHeIoklPgNI+URoW
De8yjbXAoDza8zoXb0HJCaPPFYdE97O3daq/mxdPcOMhAw6ZO7JNHKWYb5q9znZghfA+3IGytSCiPib6Fl4tLiKQWjHsomEvXb1n2bIh+i
6UPeUObmcEu+BkMv643AA0OBQkTNqyPl0LL92LdDYyv1N7pti5gLgjuj3rRvlHSoBotgS6QSDs1EA9QrEVhrJK44x/CIu7YAIrjwCPoALWM
GHSK4JBJcD7N8DDeBY+r7YNLdWEK6TLqztEeiKFSYe0n2AAqC/HMdYQiU9jIPilg6BdaNOhH4sQOAZa8av0x+H24R/KzPkRQrUr/RV26E
QUSBPj4rrPaOLV3Fe8GFwE3rMot4yY53gZopxFLjawSf6fKglQTK01yxNy+1zv2H1iATLqbHwUcabS685qBDq5ggz1fDaiGy4NiASBX713H
mv532A7ary+fwhK5BcpGvjCp0zJ/idEsFWKN5Bidteyl3o06uJNYT//plPJEkSOBJmRI8xY+6JWxiejbGuA2JhXSV0jIUcTvmaiwS3mxed9Sr
k2HUp+TtJLKvywucKiKir+NjR0if3xaKlCR4mseTfM17e+WS4XpqBPxPwkZrjihxHqLNj7qGLC0FpS+sccSjyJvHWbFGBfoLkKM+EzJkmBEY
DVIA5C/avk7f3kGnzqL/gLhyMMkdgJdNQt3naBiawxyrBhAc924pIFRpI8qm0XZUULy9bGplVEpGL7t5QFWKqw21r+UOFAMfTWVDXg3q
Jhteh5yBwhiND3udZLPTx/Q+wSJoQYlkSqIkvCBaRSBn+BMstzNu/NKrChMcuQnZfYFpaurfmY7mQFHQC7vAexl5ERL1r8r8UGCxlP6CJF
4cjI6hVaR86zIycDP/MdiCmZIM80CtypXIzMdbV6o1CEXE/QzEPNCH2kzxOO3x7p3+GrMRB/yX6Wyi/1kSxr2ePx2WV0DlY6u5B6HOhxAL
v9ZEbZDo2Y3ytX+YPos/H4tf6fib+I5i+I97WogAz+Y5RhUqADabmkRKJrRU3J60g8uNidHUJJ7sTe+YjU2BDf5B3iZGBiN0ppsSVzRjIEheEZ
2RwOxvxOnRFH95vKRWZLAjOxGXMyY++yfgqpUSg4CFUcczqmOO03ufvmYBj6MC4L6dDrX93wPgswAeclhyAyXA5f8/UgS8pHXFNA6N
8DtQrmkLfzCjnIUtouLvus1xItk2DG1D8fxtHV0kTeJmfNrMgdsHQZxy1XE25PQme+UHfPCtFrT0R69hkB95rBX2lMAj/juipJ/fWv830XxUqb
EFZEGlcxJgHOujYgxmbQR0dX2LE8BN0Kgs9y+K9iEU64NheAHRul4wo4RpSrBaMTWucsVvyzM5TdBT5Tfhb1yUc8rtm8mtrwhJ/PYYCZ
Cac4c+qIRR3hnRr2aDDoQejFuae8UYh4S+NrjVK4VtBlVzbcX0+qf6peEppEAjco1qNWOGsSCYtZFBUqI1T682/xYDs1eBmS0QO2Q2WGKm
z/ThUW6Lir8KOhY5UBBO3NNuucdfZ47d+PxIb1Eqg2iFmITyJSjksuLl7Z3F3MnILuWeJdjovQlRYn+p0D/S/V9R+Z797JmCEh/ADxmVirHt3
oAFX6XLlzE7w1TJRy5T/zFd1d6I3Rz2+2cBjDRjq0kqOuxfSj14eL1ZJCXt7SqyX0MZnZSqums4S9//olXIQaevgX7bDe2zjpU18M93IdLHcF
PZKqX3cR7MTHAom+xY+iwicGwLcgLpAFEdNYop9dY5skaMFQTcTA3WKDzNc8WMNqah8pwQOO//ciafOoSAFA7m9flpJGyt+T8GvNK
MTbZE9hcSaQEIyrnjBCwAnyHzzawJpUj9gNHSUxMqCAQqCBzjQl5G/Wbuy+fMQ/Nf1c99Qh5tUYstscl49Na3rhcLL1Yv2NT4FziVG2GhA
cA+3PHFJueC/kxA6TSIQ7+LE+8vaMxEozdaRKqWeiMX0eyY/rfr7VAOwPn0t+k0cYhf2bwBAUmwRFKpAGjUoyiXmQfg82HnzNvdf/PBRO
Hf8FDcfF55dn0+9FXuL0bKkiiZKrpSQEXUaXtlmeucAEvcW1zz6je+uzcDfmBdMw9NscsDzJRVZPFW8D/pgAjg6R/drInIXAP5rOkZVKkj8Q
WTFCqdp/C4WKbwmIAFJsM4ljQvs4JyBuYbcCJA1eefC+fBZQwgryVgZL7wZDi3vUJaaYTLXLq5kalL2OFxMu9mngDwtc0wHHus9GGuJ
59zIKsy//ZWS6A3y0r/ZGlAcNYIXI9twvlf7ZDG97xTIFWncO9/P0ATasoT/c/2eWpqa9M2xiVkEpIipep9R2o87WcTRoxYxfqN3qq8Nh8khCt
wIIhTK0X18MX2XuZNVUk+uQmNoXrSxUQZhIMitjPE2XtMB3kMLGdnctOfymdPSPpuOtIf7zG6I77W1+BFHLVxceN854JefJEsGGE5AAQ
yIQb2jX3wtctsQP0l1jJuiMECIQ+VVzqxlO+yLR53B9ccU5JO8CUqJ2o1tVCpbpc0MP0UfHLS2k3kXotVE2sPGORZT7CQf74lifE16KeehEeb
9v8WeKi9R+zCjn1clUEcDzJbb7S/n4LPgTjLK+cs+FvnSVPr4l3boIltn8VFxPw7R906r/IdApbmYiRXHMAWzXoYCTagpJpM+q0KhyIBiChaz1
N5b1BT7o1kuzahTOY2upyUcoWM5XbdadlyDYnVreTxBtZ48WunU1zCQXTAC+P6s5ZyCTADbLV2YGOppl3DtXSKesZUHM8X785e+ST17
3MMzK/erSP7e4YdGGANFudKLOjA+5xim4maAZPXQlFxeOvES10Q8YElWZPCx/yzWSmq2P2DaStSf9HW5QtmQmxSO6Tz1tZwPAPbSlg
JJz9stmvm5xS5ppfFBcvZkLJh+5anYc6zm5mqyLR9GOAZIHZ1u9vrT4eBBgKG0GhjGGW3KJoC2X0DqqNRLsYLQ5ZLRZ3qSqyv0lzcRW0

0R+ildCUPnAPh8WYT6Sevp1WS5mvLVjFtYCCOFiM6Jz58fHe4ytflZanN3pTNuY9WaVhFway5uQ6ahH2sYAT2UmHUSxoI2OCsYdokXY0y
69JvBfs1cCVTposqgZq+Rz9UsCebNEBR1KIK7iqUzaohI/A395EXNleecJeHE74EYdn1WShkZP79w6fpNohN2EECUpckAmPl7SB3ne/+NC
yPVNxgh+Zo6MuI3NFb1xf09Hvc44RbvKl3boiIM3+DV7bnd/tBW8jm0/9KhBJws0oCac3nr0N3BWW4Q2NPbGityRGrmgVpoaIc4CR8sDX
iXlFCuw3FezCZki9/vJtNzbV0eRNhFhxD+p6AhoD9K1UfOP7gBKmz529txwidecFxjiCmBKux+GQocViz75geqqXkGVxJZ69/+O+GkOsu8C
XjyCf9R9FLeZbnEJBP/ITivd1+uGWzlQULA+PXL1edihnNxnn14yOaDCLQ4ywJTqRGjZpv1NKp3VahPVeWJG4S/3y8qe09dXk2e827Pk72
Lzkpa69NHSOaaTvMXsBbVFij5G4ELbwaK7QThn08QmbI4KuvYwHGTR/jIl5OskMvzXInLYoHsrvs8iSv61dIG2T7N3laCNyvHPFZxJCtuGQ
3nw53yPbZDMjlfqeIVFN2HjBm8SHWuXbk0VO4BLEG5Qz+BHu0oNWDcdvMLCdhZdoZ7M+EXoGXvEAdczFJgR7OzHmsrrlaFx+S407EK
68dif/lpQOp3H/I5o8sJ+tIkL26sBX4R2HRkHzvEICT6uFOW3WZWvlZcsC03AxpYomu0lx3wBuqVU6WGG21wqav2vcRUMR3fq+PDOCZz+
rOoQytQaXScHfVA5ovFy0NpJrLXpLzeKwS/qxcuygi+stIgXk2kqIQFFI3YXRa4YXIXHGRw/11s7w8M/mf8ku9O3nJOIIGcowRvqCLADo1t
21Qfrf10E/uzet8U7o9H7snhp6d+QNywTMRjHbtfkB+5sTrMB0frlDMBlrhUx8N56KDVwJ0okyhhPQaZx4zeTbfyOC+8iiQjlck/G5QRfAcepP
8HA01kHR3FtCX8DdKb564agn02Ea2sV6JDFkz86lKwikHGNaexJFnCjy2ffrATd0xA4hbgeqDsiB+nLVxmdZFDMAcAQIcveJT31wassJreQ
ER66wHxYo3idERkGFVFvksSZmS2g3hzswhMbe4KJwny4JK9ihhTJEMr6UAgWFinlfemlKR28uaswlljv43u16G5SdRo7WgCpjXysEkwHm
aiC32qhqChGBKBHYircKrOY3L3Ms8MsiMYKilXlwPWZZ7iC+NqmHEQxalo229hyma9TTkwmWrxRkfQhUboH/gZPxMpE87C5AHkMK31/
0Bw7UDuar/glxqgQitmaCn8i1PRdGCUUcN5DV3hPhPrDD3cy8SC9gIIkKQLrW1W4ykLixdQPva1kx8EHuGinSiVpfX1b+471V9K3ODSEC6I
/3Jva6DeAmt2VRpDi/Mu5V4uMofbEWX3tkGt4KcsY8eEXdlvzf7SQWVFj6fmgzhPywFSqnN9G0HN9oce7N+qBGMMFQpHUTitzHS3vuift
6PhqXiXuYyVxlaPM0BtC7vwtmteGloRyMzbu/DwimU1NDYG2TQqMr3GwLfx3H/noTu/f1k9FQajVVw4bDPQ+3UttGIVdzRcpfZUq9bsvw
GOzTnTJTZMCr4uYVmnh0pwUENu8+wyXqTj8z+p9yD7gySrNTuNz18ylMrYhUbkICmRXbIJMAR2DuN/Zo8BJW3BBSajV37zNmCfi5tubj
q8C4g8nXf2Wxf4B94b3lFbFt9G4j0+7ZN1aCH3iDJlLBm2PLeQr+BElQbEXi5ksvxK21yVl+hniNicmUIQNBQcI5vLimow5KuxqXtvzB8Td7sv
f2dfyfwkLORc9u7/Yv8UVRof3pHcDkiZgmT9UrzaWZ5xTrlrQQ7HnkxKjzcCzD6VAtZiQzbVoyZoCjHBLodLKLkX7YBz9muB7ziBMcHxgCvB
S+es67Zsgok3DP+JUwGUY9FUIAszTQuhHeruH5l7MPQlyBL/pjLL0Aq5nu0a2mp2Pdk0ZwpcFyW5o419tZ9Cg9nCfCzSLNIuYP5Q3fRYW
UbLfaWKQgEaeLhuNH/zG95oPORHG8+dPehOKU4Leq3ixYvHVeQzwAbgtSb5kytw4/Pfy3K4Ye5fFMd7x6dO4tcEWwQyKf1dspWDJScm
wB9zSojZd2rgNXTEMUXL/4q5+fC175sk5eiSy15sNpcVYSi/pkQPAMJCsg+YYqtAXnu/ABMpzWo7FucFjdou6R9Axdvrpr7X72X6J830CpaE
mg7P++LySTU4D6ZhoS8DMz5BDvFSbHU9TCVYrV9mRtW2PTOdhI2Fi+8nL/6yFN0GwHa8F76Y0xok/BdBULF8m5hCRdI+wgN0cl0lGJc
zAy6bRSEYMOl/wBENFOVYrj+C2WYvS0w5G6lVJbUstoJLQVoCBtPLz+PHUyIRhYlPJxqbIRd3Dyb3L6nBjOE+/dmvo3SziOe8WObKMO3
w+5taoLUtGT9qdsGtt3w+j7D/dDLfZ+esHBGNmGKucjZaXcFSbCqXtGxmii6qNxo8Y1P/WE40Ad5r6aSwAs6Pi9JEe33umggAe5eFQfL59
dA/3JeDTrU4h05y3EeqUpYbwpNVxz/nPv3SRimzMMs1n7yZ2/CxnJ02e8d/G4uaDvNd7wB2CZsm1HlMzOV8oRYD9uUM8QGu4/7A4dD
AyRqqJPBHxAeJsRM5G+mvuHXUa/GQl8rdi+zXQ6xEa5QXs2XPbHiOT0NcZGs14yuEURQY5Gb0WLXKYv5JcJDH+fJGMm3CryMDY8JXlv
KYgs7QSnCNXYDMXp8322Iyq0PwF4UaDGM5YrSzj2rRavfBGn550zQ+wD5XpHj1NOZZYCcgBSlEEoz+wDC7KpaVlolAbFVEP+naCP8gLN
k8JTfNXMzIDlAx7S+4jr0KTBp6gWmK4Zq5T0d8yXBxYkNmpyLXGBltnxnUBKaL+9SRb5iyiSLkIpIS0KYN/aM5cskHQ0qIDYalwqGa31GOy
w3+5TY698cAl8ICGCO4LxAgOr2R9l+1x4L5sPVvz9Bk3rZ+iWhbMQ5fyU718ut78PK+6Un8dktOcjBytb3T8byffnnE21dBabRfWlsM7zng5
2zYST1Lb28QDEjJfhfLhyWT3Nm5CkbYQleQnKjOUKF5h+bx/xRD6oycnj41P6HiNN/apkFboT1rCLOtlqFCCIIuOKaxSkTh/OgsBVdQvUbkN
5RFE8tADxpcMlMb3ZrIessp6rFYOrCuuZIC1e4UziwmSTojxhLcAvRPUrqG6rKZ4rl19mprmyL9fuKHQP4WFKZjhy8Y3axHCrQg5aYP/BO+1
4Z96CMIKq2GC3OHcbHlrsIPS3n63cG7FizJpSAduBiTqFPiI+BKYruF3PvQXlpc/0Is3JXgE/tHFS2v0uFopp4yXzFQCNiK3YtOVqCzY8DLEe
+I1PTR6m7OG9se/pdFLH6AnYc9aFKAOtbS8jmbFs2OAPd5IkXiC0cCYV6EBAQqnEN8qJO5uuSwevcaQ6agotzyw7JrXrGE462WyuhQKJ
fBJMWnLLatIY1SPqKc9eOMXUJgeI9sP9m7XrYTnMMHdi0zOqj/hPgE1Tu6XbfVwmYpMPgzhRW73k2aZK8tUQiHp8O9VzTPBVnn9HDH
KneHvxncLn0HzNNJpyWLgT7ufKlsm+Z1Ao91fYlWc5AzTHrnpbp1ZALzWc6yglp58n+7A4TqOPQr0Zz9IBYD0kjC8x7DYDflLdbJnXuW5ek
WjyFbNQ7LbTYctxSFezqjYud4uMNu58yOW3bllk41JG3Xu/WgfD1magZmHXH6JGT6i/0HCz+t7L+Uk4aCaiwN9TEA8jstJTQXcYOvIUipT3
3zIO0uld7DhRAkoT9B9kbVzcV6yQgnYoXuvcH5r4yvGwzX0svDMbOThunscRegQ5hWkY9nsDwdZpZylWAJMpOMaK6ujWaFmnpXTjWW
KSnnNja4lTDrG0b1cBPU1fXJ7uL1+eNCBg+Uh3hSSNi7uDjQDao3tu8XfcXVEfJ3mn+v4Tm+f0d2v6OSNFWlk7fPk1hQVLPf90kGdYUQHB
FyrhMMg62CFN7JDvhS2pIKV2MMEnpWKjL5Dqzvv+oTra5CylHPs/emVGJR4rFiV/I7nuS3dS8wk1ZuDzgxdlR7DdDcujZ+QestN8HnqWQ
gE1FVPo8VAZwZQj5H0/LfcOG5kD0OmLVq/TVx24F9GEGhjZhKaE0aGY28ayuhsCCJjte74tzP6E+XdmfI7xE5/1DL34gGGI810sTBWJv/E
Bkbt+mVWisFZ/Dosj75A7pFbldT1ZahmODbbw/J81SbBHyqQyIL0ueqB+CDb7uGcefleR0q8a9RXBVsqhnavd9Py8LLfHnMpMeci+a16+la
YUtEx2+nJsSg0DbuyCdwrkmKqid5WxKmlfQaENmLS9FOBxNr2eXPptPdQClA4erx8MBreBPLjjTiy3fSzx49MTrv7nIqM+/wepKIm9URcUCv
K3ZL45wTk08wBXJooK6NFriACOQeb+Mas+63ImzR7BmYFM1xKMJ2agPgjD7TUgfG4TyuDG92Ju989kfuuXk26N/3nDIfLenUS3Uumb4
31aHsgw7qM17bVcKHeTEzgeg1MjsKkl9l9fJqLqk5+NpjbkdLnAbr3pPdxgizprYd3/8sJZrmht23cjKnxGgOg+wG5AWAp0s0hGBgJOROpF
cBX0KJD2XT5WOBMKulT/GnAL8xTSQeiQYZuDW96mf7+PUrANzQVFde5Zu0NvEAO6qOG3s8Yg6Nc6uElhWMwSEkSkoWkfOHVr3tXnV
13RgUubmK1JpnzYOU1UrplztryZTRCGAnRyvhccJP0PpRPRb1vR/HJjjZafclTmLy0oMOp9helBYlVkffhIts8rirdoPDPQQDPlWElyOPbeBej/
Orl3Z7D+Q7heu/jv3H9mOVFuIJz3D8XamjOx2YcLbYM1lxvOahxbA/dzxez21rX6SGJWkdMuLKiceMtERq19mUhclqD1MBrYhUkweiVPlEf
uxgG+D7ZGYAzMvfAzNx+0IJnXbXvuiJu1kt0AbVbIGQidfUb2OM5KBlFWmQMesUouw74eQR1Qx04QeyHktm01wgueCDs05Y6n9wbOXf
f4Elm6uM0chSPjxCSOj5hISld3TZ8wQzSHGTGWWFdXOuvtIS8+Sh9ciA4IpwnPy6U0/bbDBmAE3C9NHdrXhGigILy6nfH9BwOId+NiKuJxf
8HQAu67pXpimac1N/aeJWvdianvJfTYDNqVw2V1Z4fKYMMaX0MOx2JhW9UDi7fbPdvWNdbU7aizEKdwnxImIjlOr8YAeVbnF5955TZv3s
/OJKfiABbx3BYV63UkszID9FFyCGFdQUwEz3xTQEA4wgI38RRV2avXXgUudn04uy6y+CkWzGjeDlFXmR+uM5BuuDh+6HtQK8+KagyfL0
Dp4lMCdPdJZf8iuD9FZGdlqXTO4ol7+Q7pGX1v5ysKZKf4BqNBX4BbPm7EbKyfu/lcrZP2ZuJ2WRXp7owoTVnjDSbuYc0nxMtM/BASeJy
GxogXTt2KJXJFcryJ2gFgFgvvsM/NmC/pfka2MOGOj9zB5HZ50QJapJO7bbIfsH+yqPnclpNQCeFKecoW494vdq6Al2Emc4s78dtxE5vD
7UqfC+W+qNa9x5kho1J5WnhIiDqvFLiZInU9C5r+3OeaT0ZoOvZVk4NoXmBb6bmGWkdGaaUVSRwulgNidCEW04pY4aiQVgr1rKDl5ChHtHu
Gv3mayHmy5OYcdseRtzI3YXikDYB5Lh4Vb50L1IVhfG+Nox4Z2hDXBbEDBUdCxN7h90Vvp8G5/s0Z9PUn9SaLzcaBUaLC0joJ4goLoGfM
S3Wh8ujRJeEjFzg1sCDdwBC2iZFF8/lHnc5dOpLV8GMqfeNB3FdT7qmuP/n6+zphQ2Cf4+hmYXjg2bUV/vflUw+HMqmcRWoyOtbEcEspl
EGRxiCYd2ckZ2LBSxzpube8kWwKGVpNSqTwT2a2vltCB53onmxi+orbbiwGQWAazWTD94Mhvjd95jFM5sXTN0KKqVn1X+ZHvJzpP058

GtbPxqjg1truezDePvW7qaBz0/edeDLfsgStdlF96HIc48ufG+fGoqj4VfpRXeMPYmOb0v3wW9/H+1U8VQ3IhSC5poKnAf3827XQK9WDWv
q/IM7WcHaIXBSY2NJ4q6GdtGfVhVyyUurSfP3fpFHObuFNsGW/6k6eOZuqK4wPAar1KNx45NUxlRgpT3fXJXHacIujuayAvX7nBpIpPJO/
Si1BIdjAUYqVtImXDrcxbuCdXNzBA3leGemx4G1CalxvCNW62jXI3umvwBvFz2PxA2LtYI0ZOYnwREFzapBOSgg71bzUcyBLW28/hbi5uUt
XZAKg0LdaXKBSL5yQJNon4iClvaKeGA9rwHixSuGfM1ajXuZ2jyFvX1q+KR9WV8xP9QI5IHFPDaLfn6F4sTqj6dKXm/cy0PHNwOddp8ug
dKDish7rutel/OxD7Cl8FaabVqmYhIdbsqeTpJh6RXhi7rhHKLLINLAejpI8AbtodDxcYZ3JGoMqGWLI+eGgncHMtGRBYdygQl1kRLNtR6HD
gjYsl/jv79squCQSPVZmCDux111J3rP0Gy5gFvGTkpJ+EcgV24njTmAIut8ErSk1c6gK+361tmJq/qxsKn18RFxi+p6NgUxpDomkUCa7vOTR
aNRK+VQA18d+Ozeeybf+QZXww1HFkuqcL4CtwinQbXlVQsWCdBwbwfEr67eLwEGe6gYBzrTvwCt6vOyVPQ/6O5OxyD49Fswe15cx9YL
tZIHiRMIHkqDZBisyY/U0S0et+fEGdPWeRQbNapTYdW4MNmOGS4Ch1lYfQ2D0nXZ93u0N4DGjDl4ElEx5QCuFSaT926JdKX6NSdHR/5rm
s5FnBKh4PSVDgFZK2Np3E0xgmThPv7QB6Ypo5s/D0y9S0MERXzc3MvXaYNLy9nto3vgzZ/AzVwxcMgggMHLsBZlUyPnhEa2q1/nMhvkE
K24b1R0ZL8TRkuJhAbfcCzPhpORxTulo+QFfktk8nd+ZljGnVSfS7MjSXR6dniz8ghmKH12mqtxt3V42jzz7TWsPZVbnR6R2OjtBtIDZqPj6q
LXgr53ZvTvUUen2oykR1rJQmJHyunhKkxuOQem+5VB3nKX4fLh+AppZsEO7AX8KeP1kuj7u7pnVt5qLKjLLoVvJRLaJHi4IccBCwtxXtJ
yF8vLKg5FenuOKAvtHUXu5kFOD9TksbtToOms1o0zJaWk4y5b5QNiVrTaxgq9te9kVYXVSGm+KuytANPdQvMbSRRtRnP26aq8bfdJRMv
U/m5DgrG3ge6hvOCkl8t3y2Q3nYUy6M2YqgNvLNQ62qgXCcDFUWN5vEl7PsszOTzYmaC+N2/cTCYkvWEMSgLlfwMOfRWmX2UdCtroh
oezFOoBl7Oailv9/Q35rudTiEPfx9ZJu7tiwKcLPSdUBkcl27xefgqgqayuXvY69nwcXRPkw2Hh6HZn6F8kdNWynA9aQZGMaW71OVyAZ3Y
PLovno/JGKVhJCsMQ/SEnnAEkWuO2U9eIMvWXB7y/2yFZJQcJ8SfP+vdaQXa4/sYFjjlxAgxUr+D2SW5/ZMH+8kcoAPF9E6NN3CeCYBQ
9ho0GQiH4G4obyaGvmSpGY70NLEAsm5O+ysKT+G3zVgpIsaeZk2jtIWmiVNuTuSkcc2eiuPRH+xrf0BfI93F24vfJlgiIrqsOiCtPOviTS2Hd5
Rzw5mZroIBPAa92j6hfINIDrCOlXna0iZsNPkmlFmTyS4BwxHKp/tQgJVKUN++cS8Ipiggu1PDm/yeEtEl/Lm5f5kPhDG1VkXGg0DWSTq3/
9cOZsPuypePoSWX/m0BB74g4OlluSasvCH0sdmtjZst4JzxTqaf9SCNBdl67CSC010oiOmgdwHb53geALDidqsFOjj8yel/kT59tfcwXG3mefl/
TnRIFLgTMQcr+dn4Gc7TLoMVgpHHtaN2LoNVfPtYQTlfZXSraGSXfF5SwT12vkFiUy2BOhUx5N5U2hCZJ+dt4jVgSo4nshd69mp+tA3qdZ
5vjYnjvonIfmdoz501ZgvUejYZf2te8EV2QlJ8hHkOqnEYAGCMSZjmUkgoFK2+JYmktsOgdRUtMVM6zoClBeCtUSV8uMiaDorHzIS/eQVCj0
glb9AUR8mJTUlhis5cJBa4PHbX5N/8dJ5wnOX3gi2lCTRVJNrh5eWy3wFLmXlcOQBohfRPqBa7mAraxSYLfZCuw0yxHY9vmNe8mAeA4
EPQ60+obWo9ruYvk2rAfxGxnjaQkJAj4WF4AmYTvETG4gJfYrfoRDfBt9HZhXwEsLQyJT+J+7YAXbbqtFQFBIiE0Y4kuvFOoBLj3l0F1eJ4h
M2B1WgM1CSyyf95FhzO30Fs3gMxOfOhWSObhVrbMFw1VANyd0TtxVb43mji1sNlFE3EDD9rRgl65Ov3nrh9RyF+As6Bcur/Tt6Qct+R6X
9Ok0BaE7Goqt0UlJbOno/JknL6+qhQe7VxKfO4bPJhEOF3/sY2Pr9m97yMD+83Dp8eW1pjCKNGO+FcqTpBJ3M1IxE3P4DfgBPc6zFRSu
zUGyVsKP7Ia3/KTVMbgTXx7KyQbwN6YfbfSfnSmJd7USJFGkBXex1JwoVBQuRcEMgg3qgZdOfoLBJWs8izyBHtvhGSuzHqp6W/143h1
bvvNTAZ5/QOlRAq6HFyyMTWWFZ0BA9U1xE9o3ifwGZ38N29jygMHSiI0XVtdkh3yushGRBEv6Z389dUbirIbJOI763VCvkLlYkMjAatU2iC
M4/dcohBHvsWMIOBqEbOX5YllddpnnGHPm7z15Z/ER2BVBx+dSWSBVy7xWzHeKEXD7SlqHBljCYCXiyHoiKb6Ojw5aDiG08b3qXBUC5E
bqQgSOMy9R4iAp7rNv/GaEr3w25ISIJ/0SQ+Ap++AJGrtsxBmh/EryV4nyQc60EhlSA2RB2Sw8eCb4QBKEnSAZ7M/3qeTQ08YrHZoCStD
KBXu/LgrVflvOHySQFf/oyhH+R4rsKy+hwMoYnHNoF0wfsC7T3ag32Mrdz10UB+7F7UNik56e/oXvjsXLi0mJ4XF3yHqtvjZNxSrhoZ3qAo3
jED2dG5+TeRc++Qd2dLqGxyqaLAtRu2HSdn7kXzbFSlurIDeTpTbug7WRNo0Fwo7HcejHv+dmF2WPszjn0fVdPdp13b0u5KnNLf4G5oAaZ
3NiVH7TBAlGr8Ts7QHn/ynWJmcg66JxcdZ8Yd6zGKHlJQYZz7p+TLS2daGM+Y3GWWUqWLkkrX3jDgGEaxg1nsSnQbkSWHs1Bu6VPXel
uJUEyrNQAm722m+dpBUfvySpASNWImn1+AvOWxX31yxIFM2aoGrZr7xzPmuLxABiG8VTAORItmtiB9MLFH3vBT4JQS5Ny3pSlr5fTxN5
pd4I4z6AsyNgjwNWG6eGnNX8W03KUMv4DrJR3W2amA/bY7uauXjBl5Ybe3iqx/bIqdbJ+O+ezJHnObfTearWH12/JbuUV7idI0ZrurzlHJh
cSzZqPYdIFwFA/CDSpHyzlNiqEYNiDgbTn7krTk/fYe3eYERHC4Vu2cglwH+qxFFYzaAiHDa1Hao/xOk4gCfCGr7Uho4h17I1DtovfyQ0wPm
AOTD1tvo7HMMp1neneg8rCv076v+XFFJ4PH1+2YvGmpBCy82Xp36jEXFwd5EV+g7wbejRj5IkjtGafHWw5cCSRkmOUVyV4wf3QpxxRI2
LEGx0lLR5MxAhK1nEx3aahHnwcAoO8aL44NdH/zt1vLvFJHD7l1i98Bt0vAmHMK38PgOjuXRjVQy5pbao5FhusrkC7ScLbt9XbCWYDeow
g9LS9sJQuEThod/eElq8OsmIpzYSez0OH4raX7woCdD67H0DoifXRbIe3HQ+wGACNPUfNpYiCQAh1UfKSC9TG7V/uD3UYL5HU96t3SrYV
JlQRBheDXridB84E5Wj9kYxraBZEW22S3CrdgnmrcwrLZsZ7cgDd9ufL5oDrDnKMUEHW0zDMZTdJAlKH8lLeFnnPiwe4WZ7+LyescO0SZd
aOCtpQMbViL44t34sN4+UcVm1N4z4olosaAC1pARI/4VAOwD7epKde2nJ39Ec+0aCkZS/v2wgUzQ9P2/8mToRc8+wfxBflFBmWHi1oa3i
o24Z01Mu1fblxy6BqQQ4NhZQvvtkZ39g2Bhv9IcIrZRAk+WZSstxn/CUlyn2XbnjcqTqFekA8hB8hjM1wOu1TSl2c7ra6weJpuN0W9RySoGS
3Lb8iHZSi/jZJva1YltdIBFyBRrA2GDRI+hBDhOSTCXfa3lZz11tezyld8nij5JJNqubRhKS56si1Wh6cDR8nzjpKTW37FYyBDHX3DEFmHWRq
MPrxxF+5OrBJbRNigqGKUSSxRUC3dXUH1TDGl5YNMeFmn9JgtG89aQtulPNkNRp2nm4T/5Nh0tRHlMaAcZb+gtrYMLdyHJcET3Bl+7Cf
pLWxSGAwSJfhJmcTwY1XSm8l28OSIXXKg/+wbZohai8DIcBew4rvulUk5lRJbcynkPQAU9YEFeHFg4gy/vwtwMAALwZdZaHr0ijUZjuTpfC
W0Qv0WlsXWQJh0Je9KFUNTt0nupcLAtiHXxt6SgW9z8DMJ91Mj2Yhlj/oSdl1abOzgCVQIVleun7kwaTAKeN4GBA8jgXiWsQNt4qVmFHyl
/fDFbPSkalTxEUihKZvc7e0+evdHifuPKqAjrQTgbuGAeKgs+ex7Ll9fDcflqApSR0YtLSrBlgDS+dXx6dSvZXMVt2MgU93oX8lT4+r4EV+oJrTa
+Teg8Xo87fzn0OPi3L0tt6hNlCOXLu3qAyPwdqKd9ahsAN3sbw9zmk8IxeFnCEAhXS/B61hs/XPhOFkP3Oza6DFQHPDP6YaZlWUFseOg
oqrXX8/Ctu8Hr/6Id3Mebb5YpPVxnAA+p5H7DXcGJof05NZti0DFMJDCfbzCIVAv1Ko8PKdsSHAoMOcyI/NDyuC8JoO1SofvAHiS2lzz0K4
OooILQzsR3a8OkinlajHq0ieBDJI0zs/qClGczPCdKDSU0Gj357IslEWCVXVGCMAQd+ROjg8G1F0BTa7zydwzhBdJMkrQVhQOQK9vR2ptN
NNoytkE4geSPvs9n5kvbuEkodhyGpaJxcAyQdK24vPV2dkBG33JB6avqSp4FHSdbJJoMLTt7l0/9UsazHDneKm0cJ1O/ouVNY1aTCwmjE
CxIzGIXotAKZw1e3DGBLsgw8SD70dw2iAlChQGwrQgeSsTyeNckL/zedoDvMam/t7Gr8GRwgZRTg/s/87XJtoH+i+wfX0YNoxb6rqSR/d1b
Tn5muJspir84wDyYAj1wLVAa7oWEy16EdeK9u0F6mZZ7CuoYGBO5vyo/WfMejJlTqLa+1v+l1eXX6a4oQycXjzthVCjzXuzyFjUU5XJ00DM
JS5bkmXJQ0jzSir2ax4KaYHYJ1Mb/JVkTe8XDkxVk0oDDhmEOx36+Pp39LGhvfboPKUnHVsI6/bzz0o5ALo/sBd1rH1BxS5Bgjnp9QYQ2
wGCxhWE0b46T2qH4MnMrmXQBgrEDlicJ8dMGpJblCriKvctgZOMLKY7fRU1h61XqzaQoaM2XbYe+CRnRi1tGcn8uFIQGp4i1oztDb0XB5
KJaEYIjJL7620sKapLKx+mqFfqkBB9LUkI/LDWYkCIRX7X9nrFJ16fhQaN1vGaIeHb7c+p99mtciMRukfoMWLyEjGxB91+QaTPNcT4Azr5X
b0pNPf44ZRT6KPCTfbPF8xRNpfW0bJOK2sX76lypFCIWrSTDFarqLfOMnIpCJyCPRIx9kVqZTRQRfo/L68MsLSSxa8ypL5qLA19QQs63ej
DFX4F5WdVfKMBBkcSg97Jc8DFMXgzBvXmk//FXLM2sWhJNX0cePXrKY2+n7pmre7ooRxUNSCP79KFBdihxT6xxMhzLTBzuVOhocFLk
ysmL+HlLnooR4W3Foeba9HZTXJn1rWrXaszxv1u1YyyQ4C/E6wan+Gr+BQ9Gs9WStvD/X/vTAjX+7dsSkKaHUmqaY07aodlT+2p3gX5x1
hbXNkdxAjpeUpnR/lO0cucWEPgIa7lIPT1XxHYtybuY/do8taTew/uK1DUhjy10zAjD6yojgvJmwdCH17dbbB4KQayWNNtz/EeoyvJEAvx+Lw

k/DweFMZ96w0jW4AiuizPcpP1ExbLkDPo/C5qOZdqWY/XkSez9CkW6SOR8rAC2jP4SYI/O+60r+vWcpSubAo5eTza1/RoHHhIP/vE7q2CB
D6tkzB8COJ/5X3gx4rxHvRymdkm5hhx6IAasJEQrJBw3DgYRV+mPhWnpMzGxt/5Jyu+IThv2dtGlvd35NQq5wW9ngGao1r9CPi+43obZS
I9UbjQ6e22hYPqxMI6dH8QZouOjrJqt46zavwxah9DzRVzwPhKNdD1LrlOoWUkpnygv8WeaplUCsVUTqZ1ZMk+na7wDKX8KF6oHFLlPTV
5RlwfUGZEf7y4hajJU4Vh0LlBtEOBoBxr9rBxAHswmlq/JQkUG3X+HdEI/C5/6CoEaoAoH0Djr85mb6ZbfdK8xbU7j833SSOCLZGGD1WOS
EHZMfC8AgQ2Po347ZiCvlSKer8GYRR9dOM6Uv3BoZhnJ1jPI/7HB0huMc5vUzFlRsAAAsbwKxXPS7hw/TYgsrRSdm9+XLKj4EaBLL1gG
cq9epWIIsWKvlxWkHdck0RXjFauimG0DHwye+aA7ruAQ+gvtoTBRZtP+3W4VPX7lKFO5x2qukJBGz300RrL0+Y3chUOzUcG5zpYki8b1VL
3wZ//2HTSQUB+3BpvEKOHpnyvvGZNC8eS3sMxqi4DO12zhsULIlkh5owt7LEjvIEM0UX48+X+1Dyrc1ftDaPu2+PiFIEdE4NrcspMf9wbdOI
fWQK9KpJf7ibWcO+xxBGSwkZxWM4GgX+7cIOuUEWwN+v8SmEDYuRRv9qUDn8DRwUcRcRUZTskPHq5lss91g+Cg4rQtfSkzDnkkt6B1N
bF9FXlJgPXITOdzCB1j1ddrpfd9t/8RJ4JaGiq8vPY5bRTL7LQsZDLdlz6nyjVFt5NbzEfwK/IMB3En7aeRhA9qoGJeWHcQDI+FJCyFMOpLc
gn0qSpOM2ki+YX0EI3jQRL3uYpC3nlbkblWeKnCQj5Ki+6gMW0gzJl3Ek0p/RvofW3cit103G67FJGZd2pfC9z0hi/t/MfELK41A95PZeMrdt2
lBULp+8nSLLQpjnK8Sp3HJYUYr7q3b19KUsV0hYsR0vZE6eTbKucefI3WoBGTQCFsmBC3GYqsUSPks5vOUMlf42rxQhFaMA9vErt03yOi1
zFh1/glXKJV+xk90YpeiMmuwedVbTd2EasdbibgU/lUVGZJnjTmQplPuxZPh4joA9R6aardXvo62pUvqFcsNjnyX7JN1OBjDhnpUGOBAPhU
R8ylL+Y6qybz0HGCeNhEU0ze/Nk2v1SDUHU3nxEOm2ZrDESZFEUQEHOQhB28NEZfgK3LWmDUqmSHjnN6tjv13/JmjIcTMyH/kgYzNclix
IL4nkOeurPFyPsoDCLQva/bhzlncYK/7HPjeLQ1on7iyAbAXGrhNLrYsFCGV9yx/5phQ7Q1UellcvF0nNCRUfoUakB5SyIjK4ix2/6ok028CVm
yfUOwnjGN9T3dyR4i9WDQkJwJI5C3V+unySRTGnqMRpontD+Uo0mlNvJPUDYk0v/gY5r1B0+NxkcDTquAWsnsaGMgIrkTVif4x+fhk6QX
8SJhhvlBOMxNQPb72bfHA+dh8UhSMJ/Lic5u73Z9gBQ9/4vhMHtt4tgTNa6MnZjMLD5NSQ7bFRouzvQfykwKBhe5gcGLRT/wdgb2s88y
Vk7wxZf5DVPxUQlQ0VnyZEo7vM0uhsvARamhkNo3uivxYjU1rjO0kJbhM2f4RLhut5IsP17vPj7oVUBnFNZJonVxdsNmlt09T0B1t/gZ/mlm
ztbPXwDoEZzz9QvnhFEhzRxPomCSJgrxfyk//BQDu3zXE00+LlonXjihK5SPgQYB4CQ762s0xAMiaxopK3uaQzYuRsf5VoRlgAGi3mLhYW
9Sywz4WnspzTTr0K+AUu3o4+VlX8fkdmWFVjNdkVC68BGRRJaMdE500jqZHC2rUpCh4KBUtN2IDrY6Y86/FxXI8EIcgSA2Dyoi1zyqSBcns
tdLSPScK0mFuE21lsLcQrpl1+yup7NXwrliGQahsUTTe6dLVA0w3hzsHalKvbyuEVj6HPmZyWug77TzYZhuWi21TRotvksg9hIRH71Xy4ii1f
EDWySxKm2rxvbxqLGX/IQisxFyOFggTiWYT8i1WKrViQbGasdVKnmIrFuwqalbZj/DS2gu86Bg59lJvqQRMaVZ8KLR497G0xqy4evEQ3E/w
+U7HxpItn4Ps901gUdS5i3tshvRliUZ7mAuf7tEwwsg6e1lxi20Ge8w0F1USXaB0nBbFez0I6IS/oIC6xQC8A2ykv58nIXzOB2XgBxkXrun/6W
7x6YvI8G6lh0QMWd4Thuoz2byUXDfTaellhbJ3ftKMZlkKj2kpZcq6ql5bXNiKpnZj8a89LG04XqmiFPOIwCr7u0WdHFM7oIyY/NMvfHDSOD
WJhcqVL1rOp/VU9k6bfEw3FOVYh99Wf0lU9VqCvSiBhA45RROa2+zzMeEaYBkUBaECu+cW4gTPlzlhaZLANs/4Pvy29UvHz/jdehaMvtwX
pQ8sv2yBwYdeMvgJbsYN83GhCIzYVqn//tPN63XKxog/g+Gdl3RfNtVeEXvJkFdrP6BSjsSPLO3cZ1+70EQ1vLkYIVVOlBUwnHN/9sS1FOf
Gf55JzVjJVpj3G4dSXEB9c3gZ715L3ydhLQYPTZ64QEOUC6J0Ft08t7t771PYCiSURrjeQsBS94+AkHNeoBUsYmvb7ycC93XOp12VrZ4MD
2x4gDUSKK0aKl5qaSJ6O7mNY0U7H3xpq69EjLtBocVWbExW5ZM8R+XmoydX62eYkisYSinGSmioRhQc82dSnmopm7WOLiJm0yepOPe
Rjg6V3WcrPfaXhHI2sOly0gDSAkUBMgiHD/DJIhRHU7pnLSy85v7jumXaU4/hj+GZ/u7ip4wLBObkjnz+LkLBZefzfRNvMUbzVcZyKgSvwV
HfWLlXmpB0d5jt7wj5K8AGq9zyB6zJoLY2isPecaSpbLJgCsSoJZIn+9WfH9dmH/XFmunZTPsUAO076nOvZ719QV5bG8t3akm+DqS4pQ
kLi2euKbjHGhF6/S+3Lr2d/e8hGJra1erRdehzJO8lAP9QmlH48y+iAZBMRvLeKg0Rk1rP6PQ4Hyi3bw/lB45gZTlzNn4CUnjKroWXfsEPRrpE
BM5rLn6lQOugPySO6EzfIXKQ1SN6WU7g4QuoKcuM1iAj2nootG3yoYqeG8hTjfrLiMptDbY3ABAEjB6VsgqkKnviZWobOTCoq3gp/lgAR/ta
Y0vJKaLbmfKAqoajkU5GSzp+QF1MXkz1j9coNS+CFESGTgIK8U7/gVgT3hNJg9+gH3z2KMw3auY26M1jOUUMJxwmlKi/2hncw7ZYgGg
zOx4/Vd0f942PoKfCUXC3zXFAq2RQmYJ60i4qpgpmOAue/tGlY94wQ4NIR5qrSsQp1D5mt8roZI/0OUmYNJeWaS/Cdk39C3lGLP4AKNq
d/FMwdac3iMLGJjwUM3FQbTZumSSlcpMPRJuv3uksSaI0sd9qK87CV3oaqd0hsKl2vNyf9fs/K2o2Bvn3FsJfqqN5SI5uU8e8TAxlcNtMGVt
xfRAEf5dt3XY+x97l2CNZJSd7yIH3AtHkhqQYqAuCcym7Om0saPeYKkM6PPHEezJkf3WV9NzLW1Dv8n56QCBY1dMTXHG0IHihfPM9s9
/j/o/HrS3d3QTVP3GXB60uKbau8qxnuS17KWIsCYZMz7yuixKi+ginc5hd3sl75qB4KdRCBFgUfgp7RKbhIW3cl+7TjRDSFr17S++76TGe4i+
59nXPFigvN1WmBj5L3BaQJrU+1pNsE9bep5XjAdplA27RAiVQLLbcY+6MCNq3SR/b9F4CU3jyw5G/htkC00TP8FEmibhizMZEDlzITO77ge
qq1kjczsdAqz4u+/nYBYMb42PtjDcqsS2/AsRQiso/d8HoFtj11thL6bC6T5aS7yKpsRthqfu/XoWmKyj12h6NU6kWOLQiCuisQglLbIspjqLriw
E/OwymWr++Mt3FTK7/qji4ssHIaJas1Sk9jVHpCg8PnfrhjmHTopng5M/nA8+bIuVkPWPLNvYzLPcV4SoHGkvIkH2B/TGmHFxP27X5v4e
RFKqiT3GvODignV8XzDc34pnWK+snI+7PojZPpcMKYDoYzI1yFUjhJfxLy/PkuyiPiw9QfQ7slKm3hefG87dY+y8J12lRtIqGZUeuotgeM78dV
7C2aqU+VyhgpuCdP1VMnZdKES2hDP9QqjY39H39T0WhugcUz9GebWhtA3oUe5uaUEQVUXlD/p2BeJQ3gCpugU+sOvGWf7C9F5ntBzd
OG5roQF7Ot+8/6YD6X2mYIGP5Vod2v+3MdNQW1QnO57q3pzAGQrMbBWOgrUIapSC9YXgBJRMGupJ0Ky/+zaEgn/iGQxA7gta4bi3FtA
nZcyhIg6Qbghv0pRX7IZVyt46v9qrblVTZHmtnhqG1gB0c4n+EPeKvMdSMQZ27AUwMmlMp0qzJT/bBfTn5gSkRqOB/jwDHRLckU1oI+Hll
gYiaD5GYr+6m+VWNnHX2qmLbJ2MtV2YhCqV3YGwZDdOCC8HX/knXUiL2lME6yaTgwH5NnEfTK+L4xbQV0ASpMV5ZYajr6JrzegjANgl
9fGaEjzpbbY5+z8s/QEOsA7q0TEdNtcrW+sXyKn/wWOkEK0TT+H7r+hcveVgxnAaCfmWj98Qax0SHMktuVY/zq1Vq9xxDheI2NWIPSW4b
d8EVz84FNlZsk2zTVIiWNgeLzVw44BqXM9IJNSVfWwt98ynKnYE2Cdg6ohadITi0EzCB2SQ+wZZcStrAmPUIWqhCCY7CQX1o6E8jLo6nA
E4YkFS8eFrF0Zi9sQCjau58cHQijM53mqYqtigJSlWCx+tRJqiblEpbD+8rdHYItUgPz4v6RhfW1FsDyvqBnBcwoyTR7DSAjyiGjtBS2dmJg4/fd
hBzP9AyKR1VKvWuIGdK4ksLgT47xVOmK2BHjeVSnsTVJ7nTZmfuLmzzJesHS9Gdigv1YkzCW0SEVmmlS26dLDvgVahPvEtl32mJkQ9R
uz1OPi5rAIWy5xaeFReiSs0Z6+oG9M8b12ZsBpHnutKLs87ckIS40cFa7W1ebky8TaiU49p4PTAx+/HNLNcD/dgkp9oJKNIjlIWUGBl6r9IkqU
9/GtDNBc0zDYWXOntXy7tQMybVa+1tPZ0Da1qp7qZdZbj619e2/1ZiBCrxmqUNkwFZqfDGQmTCprrDw1IkvAMUB6rQC0pSv3PsRsnkkAT
sHoGsK8h1kteWwSaBNFN2cLnqw8QuM5moUHriEult17M1WiBcLydGNxgyNwJ/TER5obOyvA6kuHNTOYcjsrToXjntI0jfQSK42quqPSrhP
ST6Zpyg/7h8dLCaHzgiO/tirWFU120yitsA7A0HBHztFf4fP7izAS5GgaN5CFKir4U+BMCBp0kuP8mXD6/xkMkrJoWCh+hkRwnbPk/P60w3
ErgdRgYltCDg5csDeNMSwNGVNsRPckSLK0K2P+ik9FOeh/7uMyh4ckjJQ/2jx+zL8P4FixQn99F5+jeK8JSJl9WtZZVB8jzo+3svg54lu8zhl+
aBI6F/ZtKybSmPRX8DHfhMNnvs6L/11XjYsHy4HkTzChqGF9SxGYSdppiPWwSUap/z1fs7OSSVhKCiv+ChEIXFmNRf0V4a3CbYyW9cIQ
XwSLLG9Vf9G6FTZNDiL2+STN62TYZj7wS4C0Wodffe+FHJng17S17oowQkyLJPhS77P0p2aV+20ZvBh6v5+EMSdCrcPu5WFcbGNn5OS
6XYLU2IzE6h5MDKvtn97naD+X353R0jGC9eD/uahbNmIWFGP8kW8d6Tc6d1XE5jiwyo1a7LosXZ2xo9exjMTpdMJQEDs7fFr8hT1ar5k+q
LaPoVyfYLvjHy7/RTsu6BPIbz8kpbH15p84b6nVp5vihR0PZkjnyoUNGrzbaRDaUzK2Veq7bLGjGlaKNoZKtxR1+JIlayCOGKWg+7/G0E9hQ
0eiW1V7ENbffffnZoXWTQVW/pCGmjLPILWEYmlEuA1GqMOuTV1AgpmuLqh5Np/SQcvx188Ki5Kj42UzERBCgrYMFwkQ2ULu1LrUTzFg

P2sjln4/cs5pM6sIWGRbCkb+r986EcYCI8aCruvWq3h6fXhfz1AilElm99V01zFTKwIUkw7KaLQFXl6t9odn6a7Woegnks2HswvW6RWAZI/Q
BXtxhMCvnq0syCc6di9B3wh7iPyPIhvdF8CtEMOvubPHAOChz7vP+GJOiWg5WTAbiNQhROGlVMMDNYapSyouA9Z/qz11z9doEq3kQBv
4yZ15N9YknlfUbJRo8rCtm76LQqBKkybRIQYkjuzv93W8aDuQEvbE4NSkS6caowMLOi+IMyPCJpZqyiG6eFXknMrR5agS9gHwlpShY5ylca
oWR9A5j00P9FSFWkJJDNgpUCuQ0nvCzD7WULDFBrqWoJ8Rv0IgintknaQ5pVGT2hKNOXdtrS3KP1YwMO5HVdZPMakKJu2FPWg3Vdc
d4kwt6NSkpwEUnQ+w1o2bEjcRHE8AhTsyxhsv/b5fVueJ1nBczbRu1VZ8oFTw/YhtDSGdyFBFhkMKtaAUbD1tglGCfTjJJNPJ7N+0yb8ayl
+ixAFNEd+yN4RW8G7ZWDSKKdjvyWeKhomOp1ZJd8aBW0RcjnAsMB3aKdfWZ1DhnyEFo6RilrWFVVuJ+L9vVXqEJR711HJ+SN0ZO5BmU
po4I7hkRk66yXlyoHE05VtMZdrf43W3dk0Ay8CkiDjP+iAMQgEmnojcHG+7Cf6D21mVINJ6wq2bzs+pgSWkaNcue45mO+ZNa7YG8ZAZ82
yLe0mNvJM38y/0LkPjXzsaCfXBA1vwxa2zGlYnnLSzbKahyAxE1Na8ft5NQ8yB0QwSkOWlctNgBd3ml5wAl4b5umKeWeHn37TlvNiLJnB
f1xk3zh5Kb9c+pYpbAXMfJwsplwJNFRf5C9ALYNxftZ7S/9z+y0O3M872CjlDvIg1i/e/qSzSqYqm3kNIJgmeheu3v4ohzFHJbS4UWjwuqJX
XYJzU4IcjumUsu+9ltKMF8SZe9XzY51+yRCbNUeHwEHEoQTTOIwOm5W+A/Jg5EUqQP4wuTnSRPKGu7qABHlDjqoCEo9qpNzLLEWiuK
UR4iyt5Nd/w94sSCeQTamFxqMU8VvKTa0Ohg2saczsCquxGPb7hy9WuDScGQH9XUpaq7p9sDgYWwuo7k+AKAPEZxzWfnMuStWEs7D
v9HK928FocFUmwXoLsq7EC+fZ0OIKPQvLeU35hMWxFA1QzmtUr0FXJdTKSL0XMVtq6hcFzSHGcK6BDy/OWejntcpSXpWXzXiF5oiBkT
VgbWULeRbnqOwQLbSO4MRp3Kkv5bxlslCWweklqMhlfzleqPVOTLa2vsdcT+beJAzM7Mls71sTixkPb+Hvs/hvUhPJVHASeC2L7a4VYjj6Dz
bw1Oh17qBKCpSOHO7UTtMFHQe0J6ZbcQ4Kgi3+mUmMyi7BRIay7NqfZorFekpGZ03tSQlV5GlgriE9wt311eU3fmIxxcG5f970pn5pNMO
mJ7PoI2X6L3vw2U4OBbcCe5WmoZC15bNnY5sC0NvGI6SJMO8e0Ga7o1VZR7ZP9/Tvn2SFLJF6Fjue4xMbE70jIbXA/FN34Rn7OH34JZ
vQJK8RhdUP/8O2Y+L8adp+IGmKL4/kl/DYF5wllFX3SRQwjmQ6bBqH8OudetgtNacrbJzqbVAh0tpmL9nOBpcbdtj5zyz5+jHYkII2WmSC5
xRWN2CL3qc02eULsyzMnzB+jlFMDDN7rPkscCUgm/Tl838HEpXuchBoZUdyUHj8sOI14WJmrjJvhD1QsKJZmd/duvOW4eaKV4KEdfOiXR
gZTnEdgmMhidHFbD9rvxgmX3n4MAgk+iE3CJ0YpWjW5EMj8sguRcY2bfWIxX/lRdokrBnkqsA7BTfzylN0T0ds+lwGGA79gmDpMipAbeXy
3KxDkYJEx+K8jx+MgzkYmFUBYo0X/oycmcqh2e9pzgU0zLZpNpGQ00HOzFnBUwNrlZKL7YXU7uaFwuxN22I1BqDDgWt6w00Po9sN0Cp
P1zzLDfBK3C11oERBmKt3X+d4e4Sh+Xn6J4aBp8HJaPQRdpQ7JNCfFO7O+yhWuu86XhgPYkNybD6a+QrrxuPxvGirZiWfTyzE9R2zPBU
wjKReyQKbsBCHQEmJkTb0KuHy/U0xbHwgubjpdHDWzb2xR0wv5oz6ypFCvDYoe8Ujc6pG+HXl0Urox8B34S+msCKiOPXPm0LoLHULW
3FWiloRSaYfne0YNmjIqRhfYHGvWTZdOumqd7M6TgYvcAOhpdzyT5Jiw/3DL+FiD/4DJvk9zD6PlQfIBifhpVs7MhJQ3FCCtek6MnFiGEtgL
BbWDvWGwdwHyTur6lYLAptewK4F1/FXQjK6K+JEL7wLm38ph8NhPTASTyhtUQA5SHwVTcCWHyckvxllFzieW66RDe54S0Xpp74HCKca
Tmr6c5gQs3QvQQ6TrUe4v+YVrlwn+aCM/u+cEExggiXuw8MkhAyidwcz6sfNcBQ1+M9ZZ1lNffOJd5fBG9SYmnnMJFOugvI0j4OmItSku6
HoPsHikfg0is2+23KC2aJ2SMGZIiQNDyGcNiRlXl5f+uDB3ibHiNvoEkJAbNphZ0a6mi2hgCd/+OQp3gd5emylopagIkuUjQww3ZV4yWYt4fir
8hLSNB/B0DcU5BXvIbSAys79U8IfZQShwl/bFPlm5YhxpsbnwYzyXRfJfRoubybxA8rxkStTzS2hR2XJQIKBmirZdlfwDD7Y2jFsF3SbGHT8v
W4AUy3tIbk2NbPJLX6e+fcGaDRSb/P9rqNti8dGaMm9sPMssWkhAVF0gHn+jd80Re7ie+JVwLqo5sCZZudpDpXGnzkfJmmcXaC1kQR6
mmEzqcijUa2/jBFDlF8mvJrUZ0C2gxtKzNmlellK+iNRdRgKgEHBhJhdAfWKv8CWPrZ7HbY/GycxiLJ7lPrAG8q7ao534GSH+yesvYP9IPjU0
EnaZ932dmQHQIPUM0Ae1mKQ9nM/x+B8T6eepZSorMwg4U8t/USXHt+zl5KS9PAXJ+J3bS4xRXVgYCiPMnvv5K497HguJEJCLzxFpVcL
nbaz0xlVYNlYvXE4D3I8p4Y/x7ucsnV97xYN8wlpIPHr56bsKgAtSD93ENn40iBtcCQX9Pnmk/XLbDRz7gkDFhllhKK4EQwiVf8WWvOkS69s
xAnG3vlsRa3dLD0vvWgk0LHU855kSsTy3gwl19bvrP7JQcTQDnO6iWVktVyQqCaQk1+uZ3EvZVnCVfvEZpmkE38Z2i2QFNSvR25YnDGt3
Z3TMP3T2peBA5bfuSqtI7Uh+EC/G8oXStq06vY/BNvxuCQpOuZ3t/lv8MytdZms4y2iFJyZ5oH1Jv7+P9EYETYaYGCllU2dWFb83BzH5qPJjzv
EWxMWsGdPuCD9h8L0EtodMSjUsLoYIQTNXjH6/w1lHArZ2igACRaz8OdrDDWYSKdJiLMr+B/ChC3liFt///n19W82oyf6GZM862YUQiBpW
yG5iyJ350wYpadl3jzNGoEG616zPQQ46cjh7Sn1oRc5QzPY3/uxp1zT85t28X5pk18WBn/MhQnrfwey+jDdXe1tEdPXBpzm/BUu5pEl2SJm
8mgngSV7OithxOZynnsFTsuKGy6HuIrpLhuwLkG2FVRp14t/ZZaCKT+ZcijI46IFJgRYkplFiA80IFZz7t18Z9y7CbSYsRu+I1FQMgC6Uz3Y64
9sGTxeNv1OsDS1PacLhYA1zK2n1KoxlMCYERODJbT0P36BEYhYuf3jr6FIEBLyAZZUQTordK+DPYdYAev7u5Z9W3T9DooMlR3JtGjkK2gL
TZKLRDAbXmUeAYd9o/ZVTjyzRNYrTElH9ajENQKY81UdNNnSEUj76v3Fe+ifDvWpwKoCeJgzPV4F3tr+o9fLu6ZaZFUMmkMWNbhVgUhu
7CW93E+cHcXGKOHjb/QLNMMOJQ/4Dy24NWheV07uCYlUMiN8aUHiAG09gvm4XBPcHeO7zxAQhIU37axigr4fTDDSKgP+/IiyC5rmHhz
1mAP7DY5vaaOmKozdSYSB+b55/NQn2VTd7uMbG1Gx1v/i1B7wjeeCKWH79blA2f72C8DLvJT3xKPtoxR0ATHRYIE6BIOjit3BzH5qPJjzv
pBV+eqY5HZmW7A0rUYRXHCxpwEyWWQyeCTuYc81umlGWWidzYnH+f69dcqeIgQeekpTMBMFcEGDKI3CRLm1EF437qmnxqoDXSKVp
smfgWtF5dhr3BUb4sAy2UymfjDwoM7vuER4m4H6dJ6iQYCyvA8jqGE20Y+ojjvkgAgI4qa+OCdKm9oxe+29pLhgxGeL7kvJ+yWd7rLo5Iak
RTLZ8TeVOec629aBHG14gvsGGkXNp21JSLxyRohiVQoBkMLB7Xha8C9SBZuGP9pZaF3f44c+98s5djzlq4UygqZx3Zj5WfOyJKnxhiLhuO
cCvwuECPooPAFh7VhxX+U/GvYVVdNe6z8woVkk5CIPIw7F8VLAYxN+7kTuJyG0q0YUf1i4eS9blgNxaC4tgRG8GzdbvZ5Y+1gKGz6wrQ5T
jUcK94X/1quczCWALNnLJuh+yhL/c/h3pAMr0d5MD6kzMmIMX2Nvss1BX4RxIpZE54UJ4zeIHFw9CVYUtvXUjvhsEWdLD0tC9d0xyvCK
W1BaDP6owwc8wx759x/5KYlqCBKBIebBJe/J5IP0ebvkZcMR6VPCT/qBxDrfmIPzx+QQFV7TqK5RuyW8BpIm+CcfYHaLrj9KYel4d8eWPV
lhKbaDjkkN7vGvBTTcUQ/PVCZ9kOhZmt4ORLcRGGZZrW8yHMFdiJ5MoMEEFfmszh2KdY1DRmZN07Ad9U7Zv6TgNEFbcXD5TA56Npq
WAwUplDk1zwS8s+TM/mOCuA79ve+Wg6kCPqExL1hndCFzyvyxmp6HpKA96+Z1to4BCYPK0wgRcZiBN2eycB2pEvFZB+JANskZNAtX9
RKDQnEF5jjnpkf0DIY8fSs/GuyKxda6VFOkl8YxBsa2mNRzYFE6MGqDPNCPewUp0XKVd6NOzmEL0CTJUC1Y+9yGcqkC/i/KOAYT74Fhx
T7YneJJVsFX/E1lxd+jOoLSqL05l7zfJwH5ieeYYfF2Q+R1vT02qErV0ma4aGV1t1XrttzJDVOFk9XsripBowDsRu2mc7uSYcmk7MiaQEyKU
YnVebauRh92sSYG7zHSRbunSeloPcK8BG8v4gQ/JNY09Y1XE/hhnXCFWbZarqFD5K1NbQqstGBRZ8yyC4RWqnoe8g9Qgk4gXiTvkfqgE
Ktl6G5nDn5UhBmg193tC+VNWZ3OTeNbt3TxXGCtC9/aRE7uWG1Ho0xn7CuJ7qakEuAH0FH1RiUtShBzOOdD4Fg6uajfLLUAlPJ7iVu9J
NrslmQtzbjhsH7R1uVkza0ukiAZ9W9cjjwpiu46JADWeIfoP4f8rypEK67GVndA5DnDQ3tJVac9Zy+/zi7IwiNRqoFlq721bNXxCZHdTyn5sQo
9c35tzk3OBjwqGApJl+28N4BlMTusskPvJbiW+BNGG3vJ3JEDNWK+GkBGP2yNp80O8j7FfidhuYW1RtlmRsRLyYy9C1BBCdvOmZ8nfjGI
WMSfl2qZqfRB7/taKD/ukjadRlaCSTlJXysDvQSr4DLe6kxu6fPCCpQ17InhjV8Vuv2F6Iz9zxgbXAdHYXnGBS/NnQ+s4dY3Fnm49jvivwIakr
PK6cvrUiSKL+20g/F4FaxyV4WG6yXfe+bW6BnQyXkcgaooRGHisQqT6G3kGJ/RI0FeT1AebsbKH2UIbjlN1ZVhdFyL93hoyr5laD3gYJXIoH
/ON0ROtXxttVKmwB4SANFOh/Vb1DS9RObmHDhAcLERcTiFpJqjrv1HIEXnhx+hHOBJgbHhMi15VXUsNPjd1wu6wbGb8tTLeG5CsV2yt7
99cj0HUKlsUJi6fN22QSSW7seA5jLJFq0HuZHdkrmpYqSaesBSs8b4nZMpOKHCSkY80ZDjTiDC5r3Bp9z/LoAFZKxUVZh6Fdnsc0xbwrDUy
xDYyzUqyubIp9vpylp6Ilx7Ry/DGcv9GMyv1FrDm3AfFJJBWSqqY3akszupsXORr5KRA+gsUTtDgqY5b0CfzUop1LPB5ODn+oSBoT3S7Bej

9jR7BGM6AeUy7MnTwYjDJScqpMKHN8XFsrrLCYdaUK3RJJVgfsqP9conLMFxRV8BC2RuDrKkP61x90tNz1gx73mOlwys7nckPb+7yiBS
QI8cntHo+gR8WZ+8UmIC1rokibnHTpzWGo1JNkfzVLzq/1PlF7FjUcmosm0SUEnm/POudG5xB5gGit3JKn5I8SG78OI3nFSTm7puiDvEO
mWSQg0v88VdqAU5b8q51ReR/E1PgAnW5+pkydzBxXCLWCWWbDfIada1u8YunFkvoqJoumW+ApAEFOYT/Zrf+U/S07aZuuWu7jSY0Wa
W6RnBF13yD2bSsRGeo7mr1MEHkGh8iEOC6DK8f4dJa5ZH/hWi1+3pseXI5xB+oyJR9uGXdGTCPUrPrYYas/S7dlPbvCTRcs0SYFnIQb+iK
MOhVIwu/IwsIm9XuYXnd8DamY4lEFZBAdT/DMvv+NoP5l2Hg91xZCOAUe5HLNZv7bJRuDZDha9h9unICwk0zNlyQufA2lnsMjbLZ01uBa
//EtBGw30E4vZEPeWm08cDNqA8uBquU6QoBMEfDU8+b98VpKW+UetmXOTWiRUkKP7s0wEx/ZtFTOriSnYvgsV2a5UGzDCAu6sk5zUM
80MHPPiK8PXRwOwOHw8Tp5gSM4PiXKS6FjVtixjNP9NyFk+NI8MyZ8zoLKplX/iq2P592ECC0E01XIM6zQT4TLkHU8UihkxvltszxROJVO
LRYatQN8fGSXD0CSb4IpUamP1A3qxAqtg/TPtY33KtLpVUa7SdKEgcEJ8qp4Mw0AMAqkZwtyTC0s+rDZ9FlL1rfCG0XUxJVL/ONvDgt3r
pH0ayuxMbUVdAm5WXcyWuBzN1gh7whsHO1g7bR+74bTmNufvX+/UTKHkKQOhUrkvhbz+ldz9vZdI5NAKBfc9YQ6/3PiMk8xsZjYW0wx
+oZMz8YpT9KxMUKVdM9tQvvy0W+EuaKApUJeRI3dlEuCDB0CfYughE3ERj0XvNP2AkKLBJRkdngrbXGgbUdu4VScS69gikk8tJFfnfBW8
0BD7LCz4gAsfxRLLHbwKe/P9xRG73fzRLUhS8HE5osKVjteg0Wgpt/idHOqFm5SbCHwn8B9zGbO3T5Hed3nT2Kc6po4vfqFfOGXmUqFt
U7KQx2XLno9P8l4JWRtBcHYayyDl3r8HWcoFyrFN9RvwXhZwWXiYV9IRIuM8ZQZIfKWGhFOc88GOm/99e2yl8asLKHrJzS83q/AbC8mG
3O3nxRWjlEEjglP0Z/pdmOfa+rwKyUOD7p0w8JB/Akw2v6VFE7Ez5vDUqCGBzFYB4E7CXY6jH48e6K6KcsWLZCZyIVWiW9P9KtG4yQqJ
m2wGLEvWaqWw9HGGlUOUQ4T+MiPK05vax3rBQL6/yUhGIC/GoagwmbUMxxc6/TR8RtxP1QABIksbKnsFjtziiw1yKBdDR9lKoJUDu3/R
fbDNx8XrzSNYRQewJukwRDLftcM7OoIxcmDQMatLd2ZrbCYVjgI4vMYfdyCTItHbZO3nEHfAf4u4tbAGLvgtrBIq5bePvSmin3hO441H3vb
xThpSeSE/lCYnC4oSrgtnevQ2sRBHqMJLv3GNaR99woWG2SeRr88+UW743FzEvla9DGtyBhkyw7KGD7O/AD3B3Xt1l22zFCTOLOW1jY0
EpVXbPCc2HC7W07ol4UwAOCziCXasaGY/0YlnhsD5fP1+WLCcg+pxh+yQhaf9z+TfVp+DAbJK4idBmrCNOUZ4cBqV++xJgoaFVoL+YEm
Mt3OHDjKIvHl4jbTsY6qmI4w7TqsEOGbIz7mYivJj/jWTl5rB0QjgU2C+Q44z2V8wb1nqPyKuFAshfhNiMk5zE4a13/R6pHU9MWXjxKx1ms
VjS1RO0QywKKZl1mOjW2bDSFfeMw/iwPq2fwTJ5ZZBtDYbkdxAKf1EJ56JkLadNHDq9cOn6GVSN1frKpKx/2oxO3xv8+yMmaNlfuKy1pS
aP7YUv/5d5fQH+MherVuDioLIT09oITUDdHmi5JC1lwVgl5ZV5Okm5FyP9dqgy/PKkOmrJIkiNbc8firhAqaBxP+3k9OuBW1GcQFQ6njzBtp
iuy0ycvAC3Imykc6GJzu3A44hj8VwbF/YBClqrio+TplVf71f0jztRa5z5L5QcyfaWTqA54g3KN3v/k6A19FOMolsiltuKYL7b4TDk8pzoDMAkV
C4/lyD01Ql5ZRiuLrjmiMpIGOrQfQf/gfZCjDgbimje4KbgKCO7ZL9cK4pl15di6uOQexqbwoL3zQv1eTonoOnHIktjlR7BhzTKVzzfujPz2Tg1an
n38ybMVFu1Gwobq9QINzHgz+Ldx5QJpny6BT8omFGqJNWD3TSQimP9bB9aEk/6xiGOavYFn22xssARwCyh43qz5GZSLQamtSfgoRjsO
B3N5kC6+odle3osrja3ebSTQA1uQ2CxJvX0tOj8NCjvdgAEBJnwmY2P5a1a0mmsoK5UgaLH6eemAHwQn4WU/cEbSXhNgcqZzuVWpkDt
K1QGeoDtydTDKa6EiDniul9wAnJEQacv9XVyUpD4RR4Pk7IsJ5pzLCNZkadIo0Tz7zXayN22AehoNljq0NS8mNkHuhwZf+F7ODv7BEaccjD
6y7J2RsoB8cmqjLfFwnag6+SeMMwqRRbZuML8rpbIN7Z+RCe0ZjMw5b4rs+I55FolUVGMd/yEJPQparvzdR+Xs2+O1IRyR+INwRbHmjO
1ugjhBq9DdbQI209KLydied54LrTGU/BRQwi/AzFUm3nDyysoMk4tz/UvyXTrrTqlQAOlJAN/uUk3MwGaMBOTX5iBCrAMW2DeGsWqGEUL
E30+DH3jJb0SiORrL5bCfJzItdxKzL+8wx9tDL2MiRRDQA4IKUJUQIOoznmXaRbvp7tr3vfzyUM

SAFETY WARNINGS

Pay attention to the proper use of the computer.

BEFORE STARTING
What you need to prepare:
Python 3
TensorFlow
Keras

1  Prepare the dataset:
   In the process of extracting protein data from UniProt, we removed those sequences with length less than 50 or greater than 1,280 amino acids, resulting in 17,651 DNA-binding protein sequences are selected as positive samples. At the same time, we got 50,500 non-DNA-binding protein sequences as negative samplesin UniProt that are 50 to 1,280 in length. We took 500 sequences from both positive and negative samples as independent test samples, respectively. For the remaining 17, 151 positive and 50,000 reverse samples, we randomly selected 85% of them as training sets and the remaining 15% as test sets to participate in model training.

2  Build model:
   The deep learning model is composed of four parts: coding layer, embedding layer, convolution layer and Bi-LSTM layer. The coding layer represents each amino acid as a particular number. The embedding layer translates amino acid sequences into continuous vectors. The convolution layer consists of two convolutions and two maximal pooling operations. The mission of the Bi-LSTM layer is to grasp the context features of amino acid sequences.We use the Keras platform to build this model.

3   Model training:
    The data is trained in the built model, and this process is carried out on the GPU. At the end of this process, we get a DNA binding protein predictor.