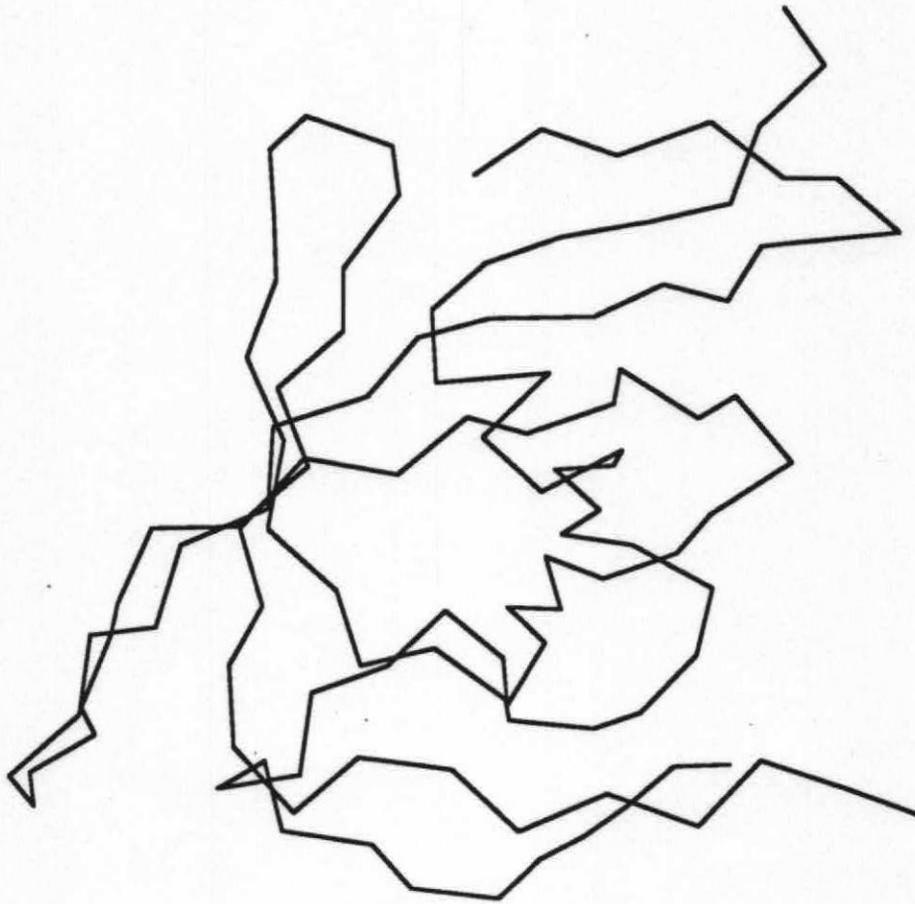


PROTEINS

SEQUENCE • STRUCTURE

FUNCTION



The biological activity of a protein can be understood in terms of its 3-dimensional structure

Structures

- Experimental determination is difficult
- 400 are known

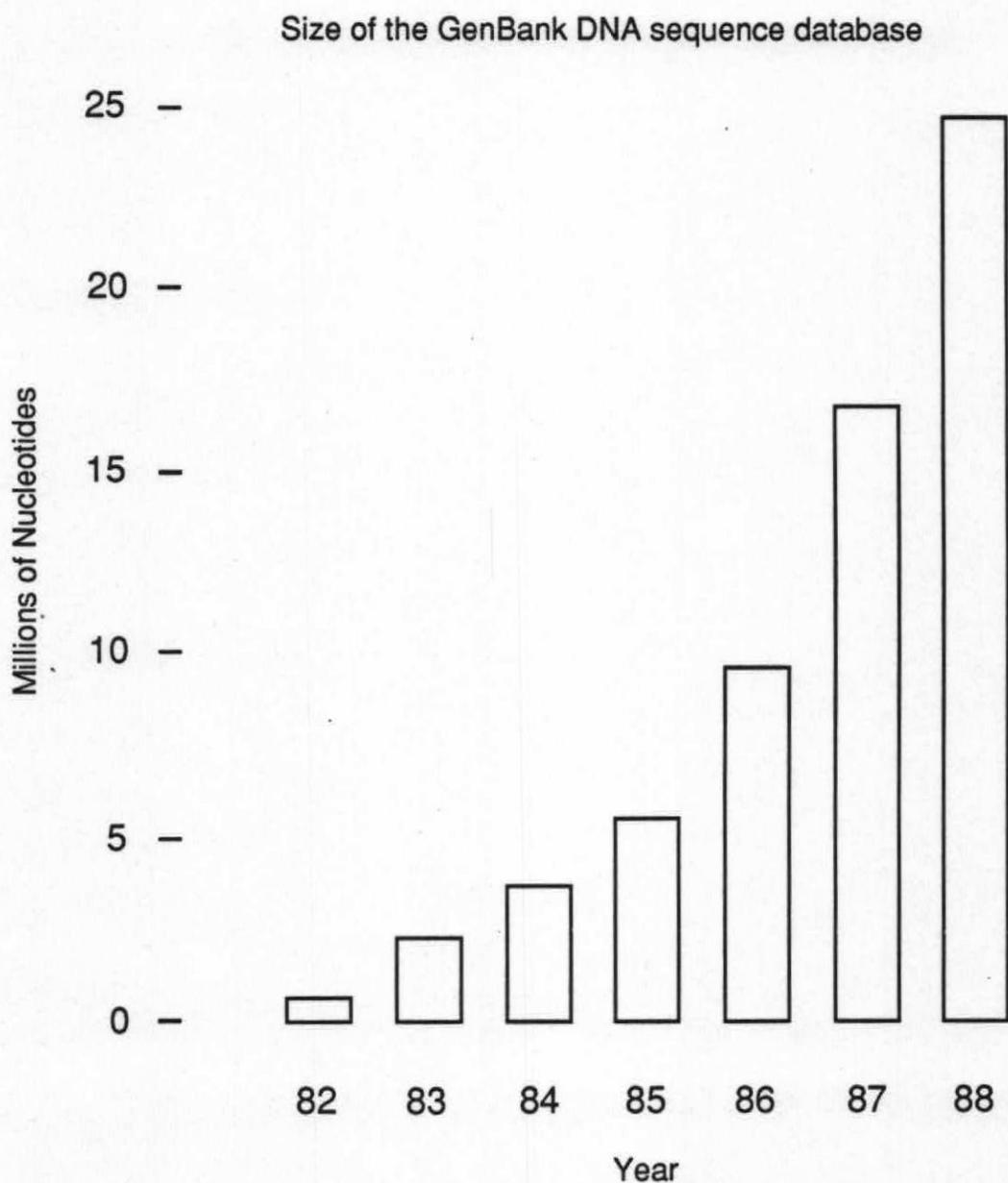
Sequences

- Determination straightforward
- 10,000 are known
- Contain all the information required to direct the folding of the protein

The Major Goal

Understanding the rules of protein folding would allow us to predict structure from sequence

The rate at which sequences are being determined is growing rapidly



Finding Similarities in Sequences

Motivation

- Conservation of a region of sequence between proteins implies functional importance

Method

- Dynamic programming is a sensitive method to identify conserved sequences
- With the CM we can apply the method to ALL pairs of proteins - over 50 million comparisons

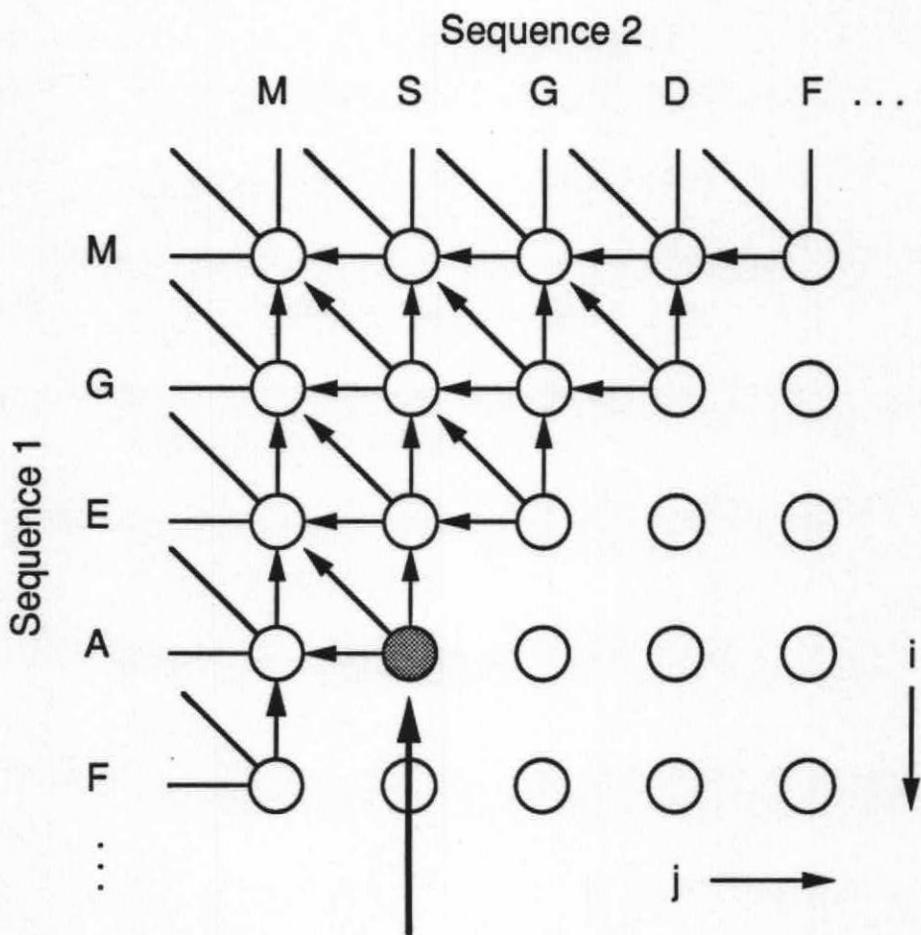
Application

- Families of related proteins can be identified
- Given the structure of one member of a family, we can reliably predict the structure of closely related proteins

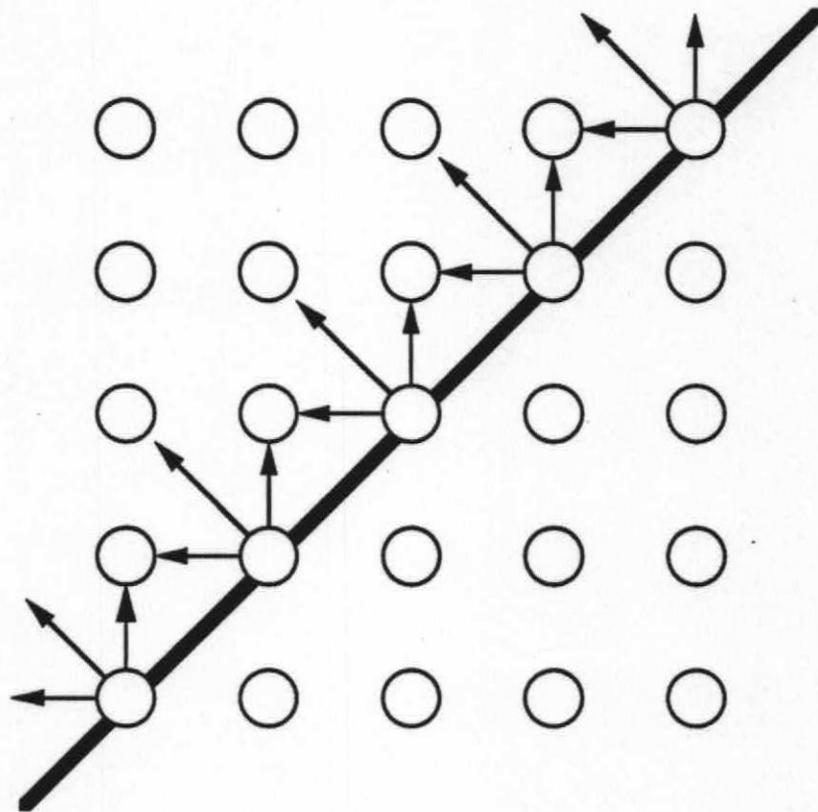
The recurrence relation

$$H_{i,j} = \max \begin{cases} 0 \\ H_{i-1,j-1} + s(a_i, b_j) \\ H_{i-1,j} + w \\ H_{i,j-1} + w \end{cases}$$

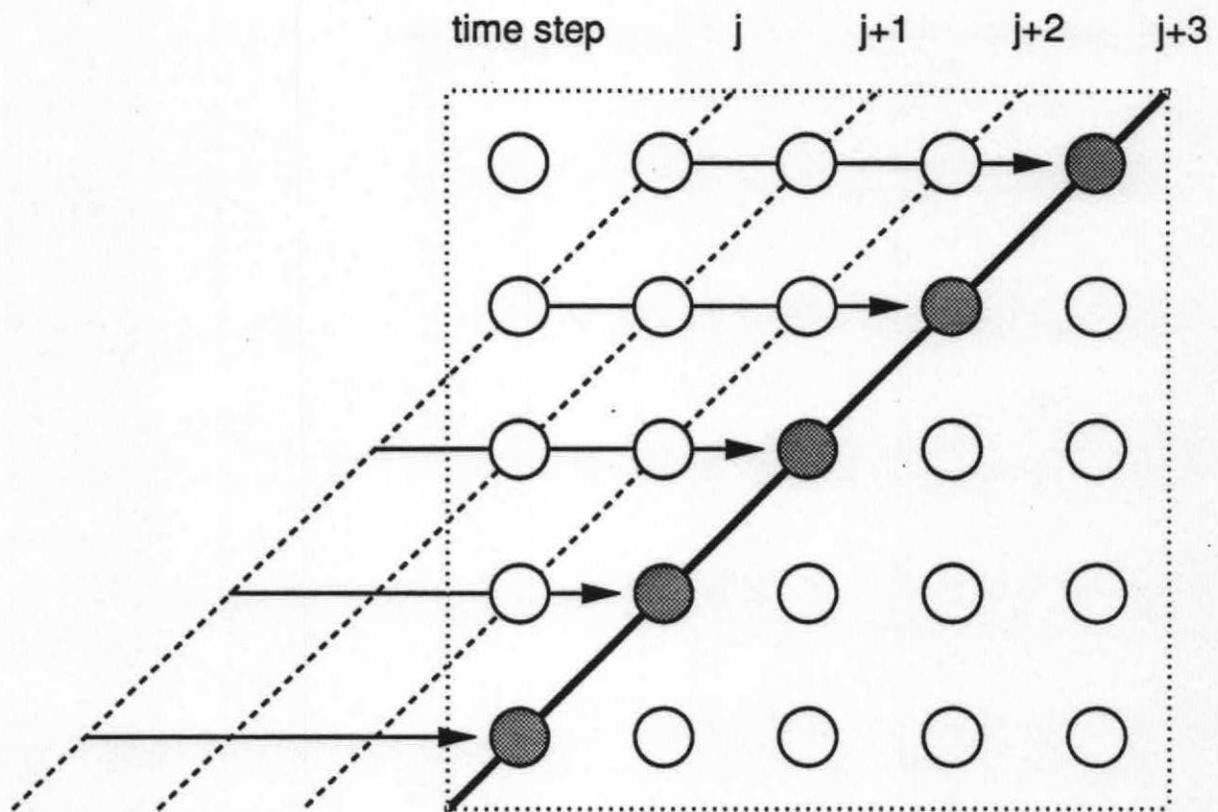
Calculation of the scoring matrix



Element (i,j) is dependent on the values of elements $(i-1,j-1)$, $(i,j-1)$ and $(i-1,j)$

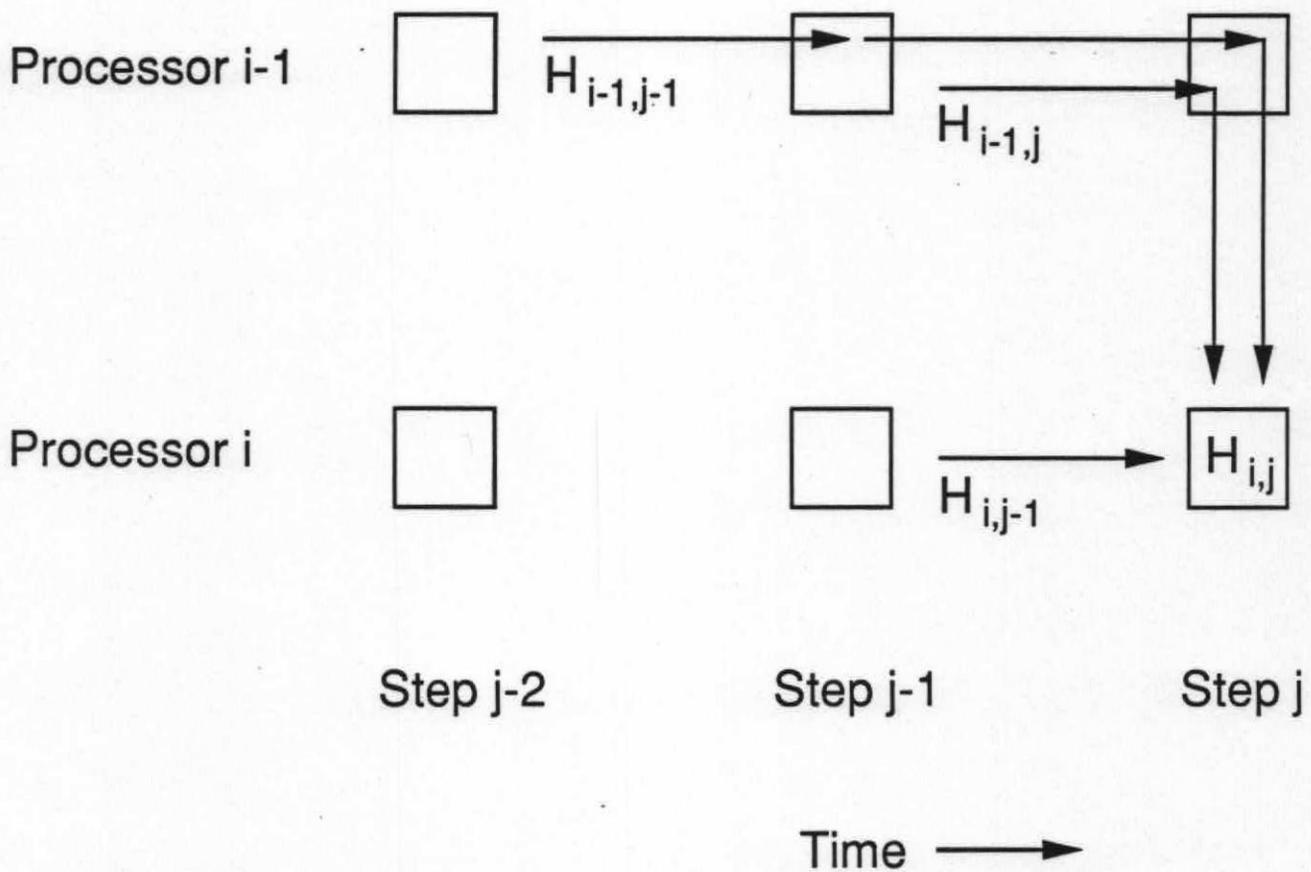


All elements on the same antidiagonal
can be computed in parallel



On each iteration the array of processors computes an entire antidiagonal; the matrix is traversed from left to right in successive iterations

The flow of information necessary to compute $H_{i,j}$



Performance

- Attained speed of 85 million matrix entries per second
- Peak of 100 million matrix entries per second

CM primitives used include :

- scans
- 1-d news transfers
- shared arrays
- indirect addressing

Sequence Patterns

Dispersed patterns of conserved amino acids often represent important functional sites

- A binding site for ATP

G - G - - G K

- A Serine Protease active site

DTSG

- A recognition site for DNA

C - - - F - - - - - L - - H - - - - - C

These can be used to search the database for other examples

Multiple Sequence Alignment

- Instead of comparing pairs of sequences we can compare alignments
- In pairwise comparison all positions are equal. Conserved positions in a multiple alignment receive additional weight
- Pairwise comparison is not sensitive to find certain relationships that we know exist
- By comparing alignments against the database we hope to identify new structural and functional domains

AZURIN

1: AECSDIAGTDMQFDDKKAIEVSKSCKQFTVNLKHTGKLPKRVNMGHNWVLTKTADMQAVEKDGIAGLDNQYLKAGDTRVLAHTKVLGGGESDSVTFDVAKLAAGDDYTFFCSPFGHGMKGLTKLVD
2: AECSDIAGNDQMDFDKKEITVSKSCKQFTVNLKHPGKLAKNVMGHNWVLTQADMQAVNDGMAAGLDNNYVKKDDARVIAHTKVIIGGETDSVTFDVSCLAAGEDYAYFCSPFGHFMKGLTKLSD
3: AQCEATIESNDAMQYNLKMVVDKSCQFTVHLKHKVGMKAKAVMGHNWVLTKEADKEGVAATDGMNAGLAQDYVYKAGDTRVIAHTKVIIGGESDSVTFDVSCLTPGEAYAYFCSPFGHWAMKGLTKLSN
4: A-CDVSIENGDMSQFNKSIIVVDKTKCEFTINLKHGKLPKAAAMGHNWVSKSDESASAVATDGMKAGLNNDYVYKAGDERVIAHTSVIIGGETDSVTFDVSCLKEGEDYAYFCSPFGHWSIMKGTIELGS
5: AECKTTIDSTDQMSFNKTAIEIDKACKTFTVELTHSGSLPKNVMGHNVLVSKQADMOPVATDGLSAGIDKNYLKEGDTRVIAHTKVIIGGESDSVTFDVSCLNAAEKYGFCSFPFGHISMKGTVTL-K
6: AECKVTVDDSTDQMSFDTKAIETIDKCKTFTVNLKHSGLPKNVMGHNWVLTQADMOPVATDGMAGIDKNYLKEGDTRVIAHTKVIIGGESDSVTFDVSCLKADGKYMFFCSFPFGHIAMKGTVTL-K
7: AECKVDVDDSTDQMSFNKTEITIDKCKTFTVNLTHSGSLPKNVMGHNWVLSKSADMAGIATDGMAGIDKDYLPKGDSTRVIAHTKVIIGGESDSVTFDVSCLTAGESYEFCSFPFGHNSMMKGAUVL-K
8: AECSDIQGNDQMDFNTNAITVDKCKQFTVNLKSHGSLPKNVMGHNWVLTQADMQGVVTDGMAAGLDKDYLPKDDSRVIAHTKVIIGGESDSVTFDVSCLKEGEQYMFCTFPFGHSALMKGLTKL-K
9: AECSDIQGNDQMDFSTNAITVDKACKTFTVNLKSHGSLPKNVMGHNWVLTQADMQGVVTDGMAAGLDKDYLPKDDSRVIAHTKVIIGGESDSVTFDVSCLKAGDAYAFFCSFPFGHSAMKGLTKL-K

1: AEC-SVDIAGTD-QMQFDDKKAIEVSKSCKQFTVNLKHTGKLPKRVNMG--HNWVLTKTADMQAVEKDGIAGLDNQYLKAGDTRVLAHTKVLGGGESDSVTF--DVAKLAAGDDYTFFCSPFGH-GALMKGLTKLVD
2: AEC-SVDIAGND-QMQFDDKKEITVSKSCKQFTVNLKHPGKLAKNVMG--HNWVLTQADMQAVNDGMAAGLDNNYVKKDDARVIAHTKVIIGGETDSVTF--DVSCLAAGEDYAYFCSPFGH-FALMKGLTKLVD
3: AQC-EATIESND-AMQYNLKMVVDKSCQFTVHLKHKVGMKAKAVMG--HNWVLTKEADKEGVAATDGMNAGLAQDYVYKAGDTRVIAHTKVIIGGESDSVTF--DVSCLTPGEAYAYFCSPFGH-WAMKGLTKLSN
4: A-C-DVSIENGD-SMQFNKSIIVVDKTKCEFTINLKHGKLPKAAAMG--HNWVSKSDESASAVATDGMKAGLNNDYVYKAGDERVIAHTSVIIGGETDSVTF--DVSCLKEGEDYAYFCSPFGH-WSIMKGTIELGS
5: AEC-KTTIDSTD-QMSFNKTAIEIDKACKTFTVELTHSGSLPKNVMG--HNVLVSKQADMOPVATDGLSAGIDKNYLKEGDTRVIAHTKVIIGGESDSVTF--DVSCLNAAEKYGFCSFPFGH-I-SMMKGTVTL-K
6: AEC-KVTVDDSTD-QMSFDTKAIETIDKCKTFTVNLKHSGLPKNVMG--HNWVLTQADMOPVATDGMAGIDKNYLKEGDTRVIAHTKVIIGGESDSVTF--DVSCLKADGKYMFFCSFPFGH-IAMKGTVTL-K
7: AEC-KVDVDDSTD-QMSFNKTEITIDKCKTFTVNLTHSGSLPKNVMG--HNWVLSKSADMAGIATDGMAGIDKDYLPKGDSTRVIAHTKVIIGGESDSVTF--DVSCLTAGESYEFCSFPFGH-NSMMKGAUVL-K
8: AEC-SVDIAGND-QMQFNTNAITVDKCKQFTVNLKSHGSLPKNVMG--HNWVLTQADMQGVVTDGMAAGLDKDYLPKDDSRVIAHTKVIIGGESDSVTF--DVSCLKEGEQYMFCTFPFGH-SALMKGLTKL-K
9: AEC-SVDIQGND-QMQFSTNAITVDKACKTFTVNLKSHGSLPKNVMG--HNWVLTQADMQGVVTDGMAAGLDKDYLPKDDSRVIAHTKVIIGGESDSVTF--DVSCLKAGDAYAFFCSFPFGH-SAMKGLTKL-K

10: --L-DVLLGGDDGSLAFIPGNFVSVA-AGEKITF-----KNNAGFPHNVVFE--DEIPAGVDASKISMAEE-----DLLN-----APGETYSVTL---SEK---G-TYTFYCA-P-HQGAGMVGKTV-N
11: --I-EVLLGSDGSLAFVPGNFSIS-AGEKITF-----KNNAGFPHNVVFE--DEIPAGVDASKISMAEE-----DLLN-----APGETYSVTL---SEK---G-TYSFYCS-P-HQGAGMVGKTV-N
12: --I-EIKLGGDDGALAFVPGSFTVA-AGEKIVF-----KNNAGFPHNVVFE--DEVPAGVDASKISMSEE-----DLLN-----APGETYAVTL---SEK---G-TYSFYCS-P-HQGAGMVGKTV-Q
13: --I-DVLLGADDGSLAFVPSFVSIS-PGEKIVF-----KNNAGFPHNVVFE--DSIPSGVDASKISMSZZ-----BLLN-----AKGETFEVAL---SNK---G-EYSFYCS-P-HQGAGMVGKTV-N
14: --I-EVLLGSDGSLAFIPNDFVA-AGEKIVF-----KNNAGFPHNVVFE--DEIPSGVDAGKISMNEE-----DLLN-----APGEVYKVN---TEK---G-SYSFYCS-P-HQGAGMVGKTV-N
15: --I-EILLGGDDGSLAFVPPNFTVA-SGEKITF-----KNNAGFPHNVVFE--DEIPSGVDSGKISMNEE-----DLLB-----APGZVYZVZL---TZK---G-SYSFYCS-P-HQGAGMVGKTV-N
16: --V-EVLLGGDDGSLAFVPGDFVA-SGEEIVF-----KNNAGFPHNVVFE--DEIPSGVDAKISMSEE-----DLLN-----APGETYKVTL---TEK---G-TYKFYCS-P-HQGAGMVGKTV-N
17: --I-EVLLGGDDGSLAFVPPNDFVA-SGEEIVF-----KNNAGFPHNVVFE--DEIPSGVDASKISMDEN-----DLLN-----AAGETYEVAL---TEA---G-TYSFYCA-P-HQGAGMVGKTV-N
18: --L-DVLLGSDGSLAFVPPNFSVPS-SGEKITF-----KNNAGFPHNVVFE--DEIPSGVDASKISMDEA-----DLLN-----APGETYAVTL---TEK---G-SYSFYCS-P-HQGAGMVGKTV-N
19: --A-EVLLGSDGSLAFVPPNFSVPS-SGEKIVF-----KNNAGFPHNVVFE--DEIPAGVDASKISMSEE-----DLLN-----APGETYAVTL---TEK---G-TYSFYCA-P-HQGAGMVGKTV-N
20: --V-EILLGGDDGSLAFVPPNFSVPS-SGEEIVF-----KNNAGFPHNVVFE--DEVPAGVDASKISMSEE-----DLLN-----APGETYSVTL---TES---G-TYKFYCS-P-HQGAGMVGKTV-N
21: --L-EVLLGSDGSLAFVPPNFSVPS-SGEEIVF-----KNNAGFPHNVVFE--DEIPAGVDASKISMAEE-----ELLN-----APGETYVVTL---DTK---G-TYSFYCS-P-HQGAGMVGKTV-N
22: --V-EVLLGSDGSLAFVPPNFSVPS-SGDTIVF-----KNNAGFPHNVVFE--DEIPSGVDAKISMAEE-----DLLN-----APGETYSVKL---DAK---G-TYKFYCS-P-HQGAGMVGQTV-N
23: --DV-TVKLGADSGALVFPNFSVPS-SVTIK-AGETVTW-----VNNAGFPHNVVFE--DEVPSGANAEALS--HE-----DYLN-----APGESYSAKF---DTA---G-TYGFYCE-P-HQGAGMKGITV-Q
24: --ETTYVKLGSDKGLLVFEPAKLTIK-PGDTVEF-----LNNKVPVPHNVVFE--ALNPAKSADLAKLSLHK-----QLLM-----SPGQSTSTTFPADAPA---G-EYTFYCE-P-HRQAGMVGKITVAG

10: -L-DVLLGGDDGSLAFIPGNFVSVAAGEKITFKNNAGFPHNVVFEDEDEIPAGVDASKISMAEEDLLNAPGETYSVTL---SEKGTYTFYCAPHQAGMVGKTV-N
11: -I-EVLLGSDGSLAFVPGNFSISAGEKITFKNNAGFPHNVVFEDEDEIPAGVDASKISMAEEDLLNAPGETYSVTL---SEKGTYSFYCSPHQAGMVGKTV-N
12: -I-EIKLGGDDGALAFVPGSFTVAAGEKIVFKNNAGFPHNVVFEDEDEVPAGVDASKISMSEEDLLNAPGETYAVTL---SEKGTYSFYCSPHQAGMVGKTV-Q
13: -I-DVLLGADDGSLAFVPSFVSISIPGKIVFKNNAGFPHNVVFEDEDSIPSGVDASKISMSZZBLLNKAGGETFEVAL---SNKGEYSFYCSPHQAGMVGKTV-N
14: -I-EVLLGGDDGSLAFIPNDFVAAGEKIVFKNNAGFPHNVVFEDEDEIPSGVDAGKISMNEEDLLNAPGEVYKVN---TEKGSYSFYCSPHQAGMVGKTV-N
15: -I-EILLGGDDGSLAFVPPNFTVAAGEKITFKNNAGFPHNVVFEDEDEIPSGVDSGKISMNEEDLLNAPGZVYZVZL---TZKGSYSFYCSPHQAGMVGKTV-N
16: -V-EVLLGGDDGSLAFVPGDFVASGEEIVFKNNAGFPHNVVFEDEDEIPSGVDAKISMSEEDLLNAPGETYKVTL---TEKGTYKFYCSPHQAGMVGKTV-N
17: -I-EVLLGGDDGSLAFVPPNDFVAAGEKIVFKNNAGFPHNVVFEDEDEIPSGVDASKISMDENEDLLNAPGETYEVAL---TEAGTYSFYCAPHQAGMVGKTV-N
18: -L-DVLLGSDGSLAFVPPNFSVPSGEEKITFKNNAGFPHNVVFEDEDEIPSGVDASKISMDEADLLNAPGETYAVTL---TEKGSYSFYCSPHQAGMVGKTV-N
19: -A-EVLLGSDGSLAFVPPNFSVPSGEEKIVFKNNAGFPHNVVFEDEDEIPAGVDASKISMSEEDLLNAPGETYAVTL---TEKGTYSFYCSPHQAGMVGKTV-N
20: -V-EILLGGDDGSLAFVPPNFSVPSGEEKITFKNNAGFPHNVVFEDEDEVPAGVDASKISMSEEDLLNAPGETYSVTL---TESGTYKFYCSPHQAGMVGKTV-N
21: -L-EVLLGSDGSLAFVPPNFSVPSGEEKIVFKNNAGFPHNVVFEDEDEIPAGVDASKISMAEEDLLNAPGETYVVTL---DTKGTYSFYCSPHQAGMVGKTV-N
22: -V-EVLLGSDGSLAFVPPNFSVPSGDTIVFKNNAGFPHNVVFEDEDEIPSGVDAKISMAEEDLLNAPGETYSVKL---DAKGTYSFYCSPHQAGMVGQTV-N
23: DV-TVKLGADSGALVFPNFSVPSVTIKAGETVTWVNNAGFPHNVVFEDEDEVPAGVDASKISMAEEDLLNAPGESYSAKF---DTAGTYGFYCEPHQAGMKGITV-Q
24: ETTYVKLGSDKGLLVFEPAKLTIKPGDTVEFLNNKVPVPHNVVFEDEDEALNPAKSADLAKLSLHKQLLSPGQSTSTTFPADAPAGEYTFYCEPHRQAGMVGKITVAG