

U W G C G 及び I D E A S 利用の手引

1988年 6月17日

国立遺伝学研究所 D D B J 利用者講習会資料

藤田 信之

411 静岡県三島市
国立遺伝学研究所
分子遺伝研究部門

電話： 0559-75-0771
E-mail: nfujita@niguts.nig.junet

U W G C G

Sequence Comparison Programs

BESTFIT	makes optimal alignments by method of Smith/Waterman
GAP	makes optimal alignments by method of Needleman/Wunsch
COMPARE-DOTPLOT	makes a dot plot by method of Maizel/Lenk
WORDSEARCH	
-SEGMENT	makes rapid sequence comparisons (for data base searching)
CONSENSUS	creates a consensus table from pre-aligned sequences
FITCONSENSUS	finds sequences similar to a known consensus sequence

Mapping

MAP	displays restriction sites and possible translations
MAPSORT	tabulates maps sorted by fragment position and size
MAPPLOT	shows restriction maps graphically
FIND	finds short sequence patterns in sequences
PEPTIDEMAP	shows patterns of proteolytic cleavage
PEPTIDESORT	tabulates data on peptides from proteolytic cleavage

Nucleic Acids Secondary Structure

STEMLOOP	finds possible stem (inverted repeat) and loop structure
COMPARE-DOTPLOT	makes a dot plot by method of Maizel/Lenk

Protein Secondary Structure

PEPPLOT	plots measures of protein secondary structure
---------	---

Pattern Recognition and Composition Analysis

TESTCODE	finds possible coding regions by plotting the "TestCode" statistic of Fickett
CODONPREFERENCE	plots the similarity between the codon choices in each reading frame and a codon frequency table
FRAMES	plots positions of rare codons and open reading frames
STATPLOT	plots the frequencies of any pattern in a DNA sequence
COMPOSITION	measures composition and di and trinucleotide frequencies
FINGERPRINT	shows labelled fragments expected for an RNA fingerprint
REPEAT	finds direct repeats
CODONFREQUENCY	tabulates codon frequencies
CORRESPOND	finds similar patterns of codon choice by comparing codon frequency tables
TERMINATOR	uses method of Trifonov to find prokaryotic terminators

Sequence Manipulation

SEQED	screen-oriented sequence editor for entering, editing and checking sequences
ASSEMBLE	joins sequences together
REFORMAT	converts a sequence file from one format to another
FROMEMBL	converts a file from EMBL format to UWGCG format

FROMGENBANK	converts a file from GenBank format to UWGCG format
FROMSTADEN	converts a file from STADEN format to UWGCG format
TOSTADEN	converts a file from UWGCG format to STADEN format

Sequence Conversion

BACKTRANSLATE	translates a peptide into a nucleotide sequence
TRANSLATE	translates a nucleotide into a peptide sequence
REVERSE	reverses and/or complements a sequence
SHUFFLE	randomizes a sequence while maintaining its composition
SIMPLIFY	makes a simplification scheme for peptide sequence comparisons

Printing and Publication

PUBLISH	arranges sequences for publication
POSTER	arranges text on figures for publication or for poster sessions
OVERPRINT	prints sequence figures using a daisy wheel printer

Database Search

STRINGS	finds files by character string (key word)
NAMES	finds files by names
SHOWFILES	writes a documented "file of file names"
FETCH	copies UWGCG data file(s) into your directory

I D E A S

Database Management and Simple Analysis Program

SEQMAN Sequence database manipulation program

get	get entries from library
scan	scan contents (identifiers and definition) of library
find	find entries by matching pattern search in definition
list	list entries by searching identifiers
xref	identify corresponding entries in two databases
pr	print sequence
rc	print reverse complement
tr	translate to amino acids
diff	difference of two sequences
freq	count residue or codon frequency
orf	locate open reading frames
ann	annotate sequence
srch	search sequence pattern in sequence or library
cvs	convert to the library (UWCGC) format
sgc	define special genetic code
tty	set output line width
help	read help files
join	join lines of two files
uc	convert to upper-case letters
lc	convert to lower-case letters
type	type and format lines
dir	directory (similar to the DCL command)
dcl	temporarily return to the DCL level
end	escape from SEQMAN

Sequence Homology Search Programs

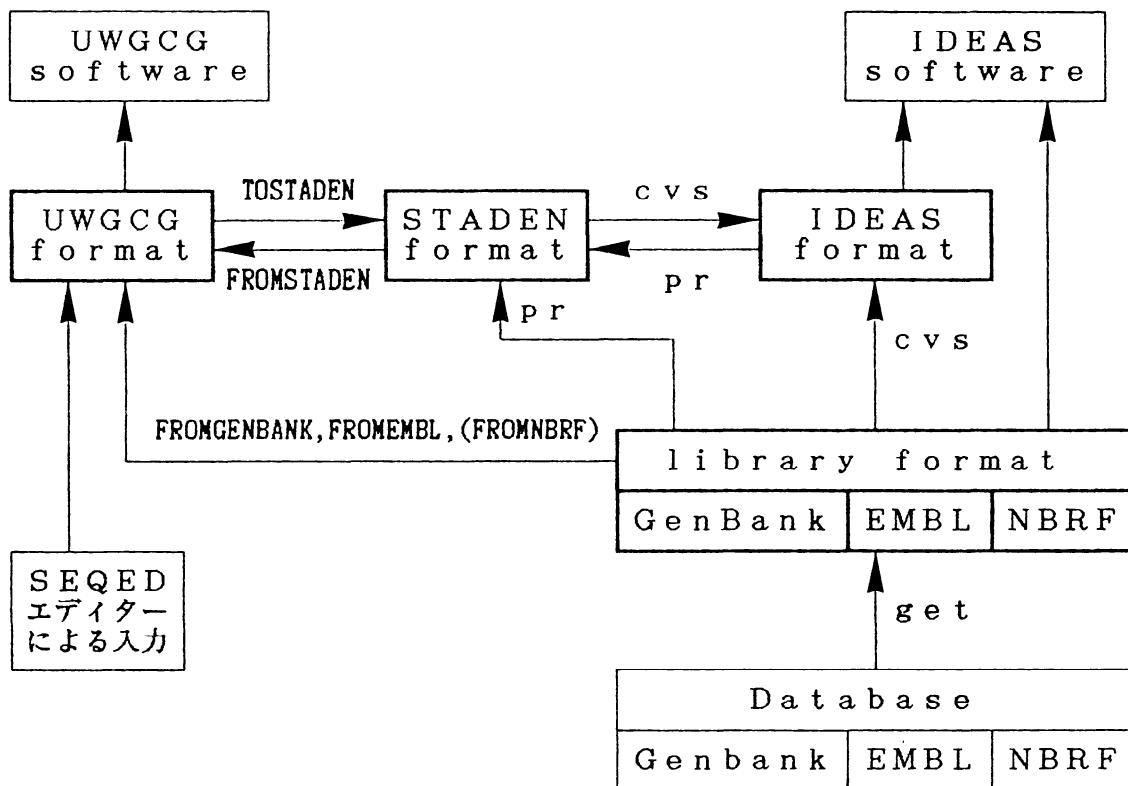
SEQH	local homology in nucleic acid sequences (complete version)
SEQF	local homology in nucleic acid sequences (rapid version)
SEQL	locally stable secondary structure in nucleic acid
SEQA	global homology alignment of two sequences
SEQHP	local homology in protein sequences (complete version)
SEQDP	significance of homology in protein sequences
SEQFP	local homology in protein sequences (rapid version)
SEQFT	local homology in translated protein sequences

Protein Structure/Function Prediction Programs

STRAL I	secondary structure prediction by homology alignment
CHOFAS	secondary structure prediction by Chou-Fasman's method
DELPHI	Garnier's program for secondary structure prediction
DSSP	Kabsch-Sander's program for secondary structure assignment
ANNOT	annotation of secondary structure in the Protein Data Bank
H PLOT	distribution of hydrophobic and charged amino acids
H COMP	comparison of hydrophobicity profiles in two proteins
WUKAB	Wu-Kabat plot of amino acid variability
ALOM	prediction of membrane proteins and transmembrane segments

フォーマット変換

個々のデータベースおよび解析ソフトウェアがそれぞれ異なったデータフォーマットを採用しているため、フォーマット変換（データの加工）をする必要があるとなり頻繁に生じます。下に変換のおおまかな流れ図を示しますが、IDEASやUWGCGのプログラムを使いこなすには、この過程を十分に理解しておくことが重要です。



上の図の中で、“get”, “cvs”, “pr”はIDEAS/SEQMANの中のコマンドを表わし、“TOSTADEN”, “FROMSTADEN”, “FROMGENBANK”, “FROMEMBL”, “FROMNBRF”はUWGCGの中のプログラムを表わす。

getコマンド (IDEAS/SEQMAN) の使い方

```

genbank
%get  embl      sequenceID >outputfile
      nbrf
      genbank
%get  embl      <listfile >outputfile
      nbrf

```

cvsコマンド (IDEAS/SEQMAN) の使い方

```
%cvs  inputfile >outputfile
```

prコマンド (IDEAS/SEQMAN) の使い方

```
%pr  inputfile >outputfile
```

例題1：GenBank DNAデータベースより必要な配列を取り出し、UWGC G上の“MAP”を用いて制限酵素の切断点を検索する。

1. データベースを検索するためNAQを呼び出す。検索はGenBank DNAデータベースを対象とする。
2. Definition行に"E.", "coli", "heat", "shock"を含むエントリーを検索する。
3. NAQから抜け出る。
4. 必要なデータを取り出すためIDEAS中のSEQMANを呼び出す。単に"IDEAS"とのみ入力するとメニューが表示されるので、その中からSEQMANを指定することも可能。
5. GETコマンドにより、GenBankデータベース中のエントリー"ECOHTPR"を取り出し、自分のディレクトリに格納する。ファイル名は"ECOHTPR."とする。
6. 確かにコピーできたかどうかをTYPEコマンドで確認する。
7. SEQMANから抜け出る。同時にIDEASからも抜け出る。
8. UWGCGの使用を宣言する。（最初に一度のみ必要）
9. GenBankフォーマットからUWGCGフォーマットに変換するためFROMGENBANKを呼び出す。
10. フォーマット変換するファイルの名前を入力する。
11. 確かにフォーマット変換できたかどうかをTYPEコマンドで確認する。

MAP is a sequence display tool that shows the sequence and its complement with the restriction map above and possible protein translations below.

Local data files: Enzyme.dat (enzyme names and recognition sites)

Command Line Switches

/SIX	shows six-base cutters when enzyme name is "
/ONCe	shows only enzymes that cut once in the chosen range
/LINEar	treats the sequence as linear (default)
/CIRcular	treats the sequence as circular
/WIDth	allows choice of number of characters per line
/MISmatch	allows mismatches in site search
/APPend	appends the enzyme data file to MAP's output

Enter enzymes: Use a "*" to get all, a "#" to get none or enter the names individually, one per line, ending the list with a blank line.

What protein translations do you want:

- a) frame 1 b) frame 2 c) frame 3
 - d) frame 4 e) frame 5 f) frame 6
 - t)hree forward frames s)ix frames o)pen frames only
 - n)o protein translation q)uit

What should I call the output file (* ecohtpr.map *) ? term[CR] 18

MAP of: ecohtpr. check: 6042 from: 1 to: 120
13-JUN-88 11:33

H
I
A
N
L
D
U
3
1

1 AAGCTTGCATTTGAACTTGTGGATAAAAATCACGGTCTGATAAAACAGTGAAATGATAACCTC
-----+-----+-----+-----+-----+-----+-----+ 60
TTCCGAACCTAACCTGAAACACCTATTTTAGGCCAGACTATTTTGCACTTACTATTTGGAC

a	K	L	A	L	N	L	W	I	K	S	R	S	D	K	T	V	N	D	N	L
b	S	L	H	*													M	I	T	S
c	A	C	I	E	L	V	D	K	I	T	V	*								

M	A	H	E
N	L	I HH	C CMT M
L	U	N HA	O LNA B
1	1	P AE	R ALQ O
		1 12	V 111 2
			//

61 GTTGCCTTAAGCTCTGGCACAGTGTGCTACCACTGAAGCCGCCAGAAGATAATCGATTG
-----+-----+-----+-----+-----+-----+ 120
CAACCAGAAATTCCGACACCCGTCATAACAAACCGATGGTGACTTCGGGGCTCTCTATAGCTAAC

a V A L K L W H S C C Y H *
b L L L S S G T V V A T T E A P E D I D *
c

nigvns\$

12. 制限酵素の切断点を求めるためMAPを呼び出す。

13. ファイル名を入力する。

14. データの始まりを指定する。

15. データの終わりを指定する。

16. 制限酵素名を入力する。 "*"を指定すると登録されている酵素すべてが対象となる。

17. アミノ酸配列への翻訳を3つのフレームについて行なう。ただしオープンフレームのみを表示する。

18. 結果を端末 (TERM) に出力する。ファイル名を指定してファイルに書き出すことも可能。

例題2：画面エディターを用いてアミノ酸配列を入力し、IDEAS上の“SEQFP”によりNBRF-PIRデータベースを対象としてホモジニ検索を行なう。

test sequence
test.seq Length: 284 13-JUN-88 19:48 Check: 7288 ..
1 MADKMQSLAL APVGGLDSYI RAANAWPLMS ADEERALAEK LHYHGDLEAA
51 KTLILSHLRF VVHIARNYAG YCLPQADLIQ EGNIGLMKAV RRPNPEVCVR
101 LVSFAVHWIK AEIHEYVLRN WRIVKVATTK AQRKLFFNLR KTKQRLCWFN
151 QDEVEMVARE LGVTISKDVRE MESRMAAQDM TFDLSSDDDS DSQPMAPVLY
201 LQDKSSNFAD GIEEDDNWEEQ AANRLTDAMQ GLDERSQDII RARWLDEDNK
251 STLQELADRY GVSAERVVRQL EKNAMKKLRA AJEA

1. UWGCGの使用を宣言する（最初に一度だけ必要）。
 2. 配列編集用の画面エディター”SEQED”を呼び出す。
(SEQEDを用いるためには、VTエミュレーターが必要)
 3. 作成するファイルの名前を入力する。
 4. コメントを入力する。コメントは5行まで入力できるが、他のフォーマットに変換する際失われることに注意。
 5. 配列を大文字で入力する。画面の範囲を越えると自動的にスクロールされる。

主な編集コマンド	
→, ←	一文字右へ, 一文字左へ
n →, n ←	n 文字右へ, n 文字左へ
>, <	50 文字右へ, 50 文字左へ
: n [CR]	n 番目の文字へ
/ 文字列 [CR]	文字列の検索
[D E L]	一文字削除
: n1, n2 d e l e t e [CR]	n1 番目から n2 番目まで削除

6. コロン ":"に続けて "c h e c k [CR]"と入力し、再入力モードに入る。
 7. 入力の誤りをさがすため、配列を再度入力する。一致しない箇所には "-"が表示されブザーが鳴るので修正する。行の間の移動は "↓"および "↑"で行なう。
 8. 確認が終わったらコロン ":"に続けて "e x i t [CR]"と入力し、エディターから抜け出す。
 9. ファイルが出来ていることを T Y P E コマンドで確認する。

TOSTADEN writes a UWGCCG sequence file into a file in the Staden/Sanger format. If the file contains a nucleic acid sequence, the ambiguity codes are translated as shown in Appendix III.

MADKMQSQLALAPVGG LDSY IRAANAWPMLSADEERAL AEKL HYHGDLEAA
KTL ILSHLRFV VH IARNYAGYGLPQADL IQEGRIGLMKA VRRFNPEVGVR
L VSFAVHWIKA E IHEYVLRNWRIVKVATTKAQRKLFFNL RKTQQLCGWFN
QDEVEMVARELGVTSKDVRM E S RMAA QDNTFDLSSDDSDSQPMAPVLY
LQDKSSNFADGIE DDNWEEQAANRL T D A M Q C L D E R S Q D I I R A R W L D E D N K
STI-OELADRYGVSAERVRVQLEKNAMKKLRAAIEA

I D E A S
Integrated Database and Extended Analysis System
for Nucleic Acids and Proteins
Dr. Minoru Kanehisa
Institute of Chemical Research
Kyoto University
0774-32-3111 ext. 2222

Enter "Menu" to see the description of programs.
Enter "Help" to access on-line documentation.

; 284 residues
TEST.SEQ
MADKMQLSAL APVGGLDSYI RAANAWPMLS ADEERALAEK LHYHGDLEAA KTLILSHLRF
VVHIAARNYAG YGLPQADLIQ ECNIGLMKAV RRFNPEVGVR LVSFAVHWIK AEIHEYVLRN
WRIVKVATIK AQRKLFFNLR KTKQRLGWFN QDEVEMVARE LGVTSKDVR MESRMAAQDM
TFLDLSSDDDS DSQPMAPVLY LQDKSSNFAD GIEDDNWEEQ AANRLTDAMQ GLDERSQDI
RARWLDEDNK STLQELADRY GVSAERVRLQI EKNANKKLRA AIEA1

```

Options available:
SEQMAN  FRAMIS  SEQF    SEQFN   SEQH    SEQFP   SEQHP   SELECT  GRAPH
SEQDP   SEQFT   SEQA    SEQL    SEQG    STRALI  CHOFAS  DELPHI  DSSP
ANNOT   H PLOT  HCOMP  WUKAB  ALIGN  ALOM
Symbol  Manual  Help   Menu   Exit

```

10. UWGCGフォーマットからSTADENフォーマット（中間
フォーマット）に変換するためTOSTADENを呼び出す。

11. 変換するファイルの名前を入力する。

12. データの始まりを指定する。

13. データの終わりを指定する。

14. 変換後のデータを格納するファイルの名前を入力する。

15. フォーマット変換できていることをTYPEコマンドで確認す
る。

16. IDEASを呼び出す。

17. IDEASの中のSEQMANを呼び出す。

18. CVSコマンドにより中間フォーマットの"TEST. SEQ"
をIDEASフォーマットに変換する。変換後のデータは
"TEST. SEQ"の名前でファイルに書き出す。

19. フォーマット変換できていることをTYPEコマンドで確認す
る。

20. SEQMANから抜け出る。

21. ホモロジー検索を行なうためIDEASの中のSEQFPを呼
び出す。

22. 検索するファイルの名前を入力する.
 23. 検索の対象となる配列集合を指定する.
 24. 検索結果の出力先を指定する.
 25. 検索のパラメーターを変更する.
 26. 最初のパラメーター (MAXD) を -70 に変更する. 2番目以降のパラメーターは初期設定のままとする.
 27. データの ID を入力する.
 28. データの始まりと終わりを指定する.
 29. 以上の設定をバッチジョブとして実行させる.
-
30. IDEAS から抜け出る.
この段階で一旦 LOGOUT するか、別の作業を行なう.
 31. バッチジョブが終わったら結果を出力させる.

