

Genomic Sequence alignments and its application for Computing Linear Structure Similarity

조환규 교수, 황미녕, 강은미, 이미경
{hgcho, mnhwang, emkang, mklee}
@pearl.pusan.ac.kr
부산대학교 공과대학
정보 컴퓨터 공학부

 Graphics Application Lab

Biology and Informatics

- **Computational Bio-problems**
 - 1 hour problem, 3 hour problem
 - 1 day problem
 - 1 week problem
 - 1 month problem, multi month problem
 - 1 year problem
- **Algorithmic Problems**
 - Formal definition
 - Formal evaluation metric
- **How to describe**

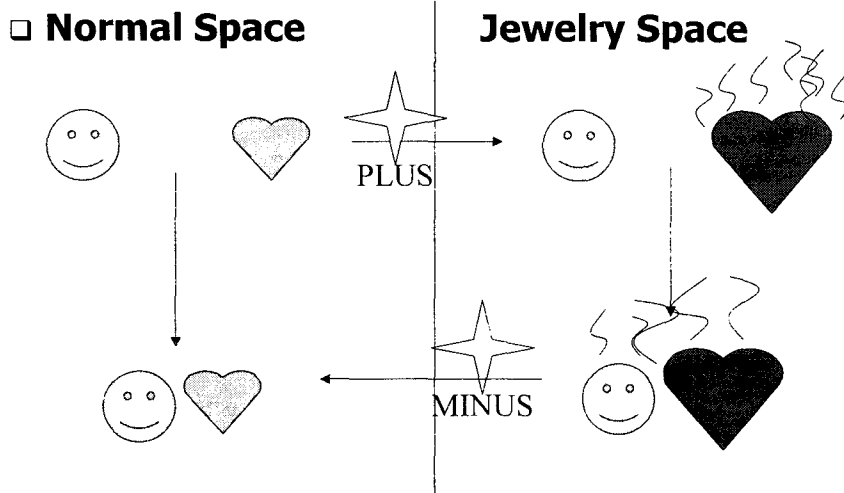


Graphics Application Lab

Main features of This Talk

- 이미 잘 정리된 **Computing method**를 어떻게 **Bioinformatics**에서 활용하는지
- **Bioinformatics**에서 잘 정리된 방법론은 **CS**쪽에서 어떻게 활용하는지
- **CS**와 **Bioinformatics**가 어떻게 연관되어있는지
- **Case Study**)
 - Genomic sequencing alignment와 program-copy detection과의 연관

Computing Space Transform




우물론

- 한 우물을 팔 것인가 ?
 - 만일 끝끝내 물이 안 나올 경우라면
- 여러 우물을 적절히 살펴가며 팔 것인가 ?
 - 이것도 저것도 아니라면 ?
 - 이미 파본 우물인지 어떻게 판단 ?
- 그렇다면 우물을 어떻게 팔 것인가 ?
- 많은 **bio-computing problem**의 기본모델은 이미 **CS theory**쪽에서 잘 정리되어 있다.

Genomic Sequence Alignments


- **Basic Assumption**
- **Why Alignment ?**
 - Similar face – similar behavior ?
 - Similar genotype – similar phenotype ?
 - Dissimilar face – dissimilar behavior ?
- **"Similarity" and "distance"**
- **Alignment Category**
 - Pair-wise alignment : Multiple alignment
 - Global alignment : Local alignment
 - Optimal alignment : Heuristic alignment


 Graphics Application Lab

Dynamic Programming

- **A Basic tool for all kinds of alignment**

- **Find a good path to "De-Jeon"**
 - If you give a good path to 조치원,
 - If you give a good path to XXX.....near to Dejeon
 - Then I will give the best way for De-jeon


 Graphics Application Lab

Dynamic Programming

- **A programming Methodolgy**
 - programming with dancing ?
 - Solution from all sub-partial solution
- 준비물
 - Objective function
 - Dynamic programming formula(recursion)
 - $F(n) = F(n-1) + F(n-2), F(0)=F(1)=1$
 - Base condition
 - Table, multi-dim array in language
- **Space complexity!**



Global Alignment(1)

- **Basic scoring:**
 - Match: 1, Mismatch: -1, Space: -2
 - **How?**
 - To find the alignment of two sequences of maximal score
- Sequence alignment problem corresponds to the *longest path problem* from the source to the sink in this *directed acyclic graph*.



Global Alignment(2)

- **CACAGTGT 와 CAGGT**

	C	A	C	A	G	T	G	T	
	0	-2	-4	-6	-8	-10	-12	-14	-16
C	-2	1	-1	-3	-5	-7	-9	-11	-13
A	-4	-1	2	0	-2	-4	-6	-8	-10
G	-6	-3	0	1	-1	-1	-3	-5	-7
G	-8	-5	-2	-1	0	0	-2	-2	-4
T	-10	-7	-4	-3	-2	-1	1	-1	-1
	C	A	C	A	G	T	G	T	
	C	A	-	-	G	-	G	T	

Graphics Application Lab

Local Alignment(1)

- An alignment between a substring of s and a substring of t
 - Each entry of (I,j) will hold the highest score of an alignment between a suffix of $s[1...i]$ and a suffix of $t[1...j]$

예) AGGTATTGA
- CCTATGGC

Graphics Application Lab

Local Alignment(2)

□ AGGTATTG 와 CTATGC

		A	G	G	T	A	T	T	A
	0	0	0	0	0	0	0	0	0
C	0	0	0	0	0	0	0	0	0
T	0	0	0	0	1	0	1	1	0
A	0	1	0	0	0	2	0	0	2
T	0	0	0	0	1	0	3	2	0
G	0	0	1	1	0	0	1	2	1
C	0	0	0	0	0	0	0	0	1
		A	G	G	T	A	T	T	A
		-	-	C	T	A	T	G	C

Semi-global Alignment(1)

- Given two sequences, check if one of them has a substring similar to the other entire sequence.
- How?
 - Find alignments ignoring the beginning and end spaces of the sequences
- Global alignment 와 비교
 - CAGCA -CTTGGATTCTCGG <-semi-global
 - ---CAGCGTGG- - - - - (score: -19)
 - CAGCACTTGGATTCTCGG <-global
 - CAGC- - - - -G -T- - - -GG (score: -12)

Semi-global Alignment(2)

□ CACAGTGT 와 CAGGT

	C	A	C	A	G	T	G	T		
	0	-2	-4	-6	-8	-10	-12	-14	-16	
C	-2	1	-1	-3	-5	-7	-9	-11	-13	
A	-4	-1	2	0	-2	-4	-6	-8	-10	
G	-6	-3	0	1	-1	-1	-3	-5	-7	
G	-8	-5	-2	-1	0	0	-2	-2	-4	
T	-10	-7	-4	-3	-2	-1	1	-1	-1	
		C	A	C	A	G	-	T	G	T
		-	-	C	A	G	G	T	-	-



General Gap Penalty

- **Definition**
 - Gap: consecutive number $k > 1$ of spaces
 - When mutations are involved, the occurrence of a gap with k spaces is more probable than the occurrence of k isolated spaces
 - $w(k)$: penalty associated with a gap with k spaces



Affine Gap Penalty Function

- **Penalty for consecutive spaces \leq isolated spaces**
- **Subadditive function**
 - $w(k_1 + k_2 + \dots + k_n) \leq w(k_1) + w(k_2) + \dots + w(k_n)$
- **Three arrays for dynamic programming**
 - $a[i,j]$ = maximum score of an alignment between $s[1..i]$ and $t[1..j]$ that ends in $s[i]$ matched with $t[j]$
 - $b[i,j]$ = maximum score of an alignment between $s[1..i]$ and $t[1..j]$ that ends in a space matched with $t[j]$
 - $c[i,j]$ = maximum score of an alignment between $s[1..i]$ and $t[1..j]$ that ends in $s[i]$ matched with a space

Heuristic Alignment

- ❑ **Main difficulties**
 - Search space, $O(n^2)$ space or $O(n^2 \log n)$ time
 - Optimality or Biologically-good Distance metric
 - Multiple alignment
- ❑ **Local search**
 - Diagonal region searching
 - Visualization., e.g., Dotlet
- ❑ **BLAST approach for long sequence**
 - Small word matching
 - And Extending from a highly matched region

Multiple Alignment

- ❑ **Problem hardness:**
 - Optimal alignment : NP-hard
 - What if more than 1000 sequence ?
- ❑ **Pairwise alignment**
- ❑ **Star Alignment**
 - Simple, but poor result
- ❑ **Tree alignment**
 - A good alternative

PART 2: Applications

DISKETTE to HAMBURG !



Graphics Application Lab



Applications : Plagiarism

- Linear Structure**
 - o Genomic sequences
 - o Plain articles
 - o Programs
 - o Human behaviors on the time-line
 - o Time-series data sets
- Student Reports Plagiarism**
- Assignment Program copying**
- Where is the original version of this one ?**
- Web searching redundancy**

Graphics Application Lab



Previous Approaches

- ❑ **Keyword frequency similarity**
 - 특정한 단어의 사용횟수 = "이순신"
 - Cosine product measure
 - Object finger printing
 - **Fixed size fingerprint**
 - **Easy to making Database**
 - **Quick searching**
 - **High false positive rate**
 - **Does not consider Structure**
 - Upon a commercial version
- ❑ **UNIX 'profile' command**



Attacking Methods and

- ❑ **Inserting some words**
- ❑ **Shuffling-invariant**
- ❑ **Easy to use in document application**
- ❑ **Hard to use in program file**
- ❑ **Recent trends**
 - Structure-oriented similarity measure
 - Greedy-Block-Removing methods...
 - Is this a basic concept of local alignment ?
- ❑ **Sample-Report-Server Building**



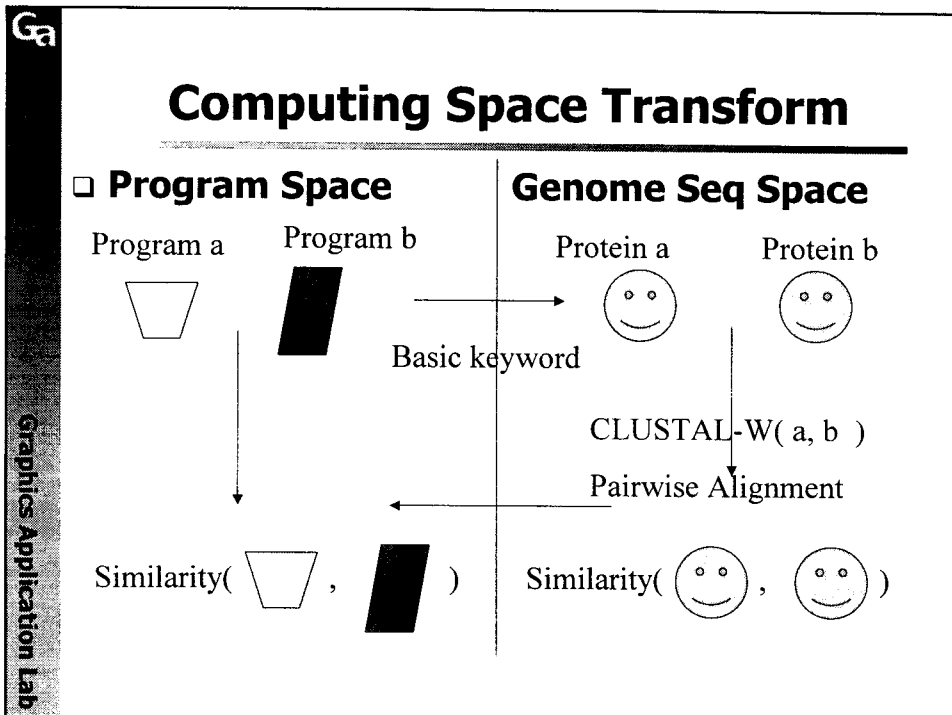
Undergraduate PGM. reports

- ❑ **Programming Assignment cheating:**
- ❑ **Assignment cheating**은 원래 일보다는 어렵다.
 - Password breaking by Mafia
- ❑ **Assignment**의 결과는 동일하다.
 - Correct program들 끼리 비교
- ❑ 주어진 시간은 비교적 짧다
- ❑ 수강생의 수는 적절히 적다.(**300** 명 이하)
- ❑ 프로그래밍 언어는 모두 동일하다.



Program Cheating Techniques

- ❑ **Complete Copying**
- ❑ **Variable exchange**
- ❑ **Garbage code insertion**
- ❑ **Function transpose**
- ❑ **Code rewriting(partially)**
- ❑ **Library code replacing**
- ❑ **Merging different codes**
- ❑ **Function resolving**
- ❑ **Function rewriting**



- PROGRAM to PROTEIN**
- **Program Language**
 - o Keyword = { int, float, class..... }
 - o Block Structure = "}", "{"
 - **Program Chromosome**
 - o Location independent code, JAVA class, C files
 - **Non-Coding region**
 - o /* this is a sample non-coding region */
 - **Promoter**
 - o Variable declaration, class definition
 - **DNA = keywords sequence**



Example

```

main( ) {
  int i, j, k ;
  .....
  for( I = 1 . I <= 100 , i++) {
    .....
    if (      ) x = y ;
    else .....
    while( ccccc ) { }
    x = 23984 ;
  } // end of for
  .....
}

```

int for if = else while =
 ↓ ↓
AGTCGCTTCGAAGCAA



Why Protein mapping ?

- **DNA sequence overlap**
 - if = AA, then = AG, * = GA, return = GG
 - AAGGA = AG + GA or AA + GG + A
 - Ambiguity resolving
- **20 Amino acid base**
 - About 20 keywords
 - 2-3 groups
 - polar, non-polar
 - hydrophobic, hydrophilic
 - Charged, uncharged



Keyword Mapping Strategy

- ❑ **Convertibility = { for , while } Easy**
- ❑ **Convertibility = { for, then } Hard**
- ❑ **Convertibility = { if, '=' } Impossible**
- ❑ **Procedure**
 - Preprocessing
 - Chromosome arrangement
 - Keyword selection
 - Protein mapping

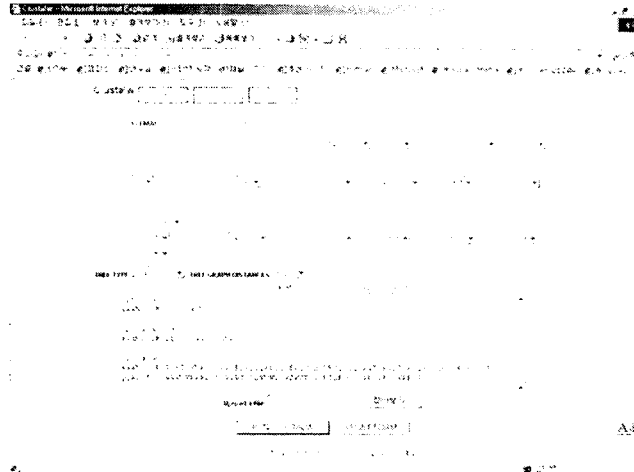


Experiment Overview

- ❑ **Sample programs <= "data structure"**
- ❑ **Students , 60**
- ❑ **Programming assignment, 12**
- ❑ **1 semester**
- ❑ **On-line evaluation system = ESPA**
 - Java-based on-line evaluation system
 - Due, 1 week
- ❑ **Do not monitor all programs**

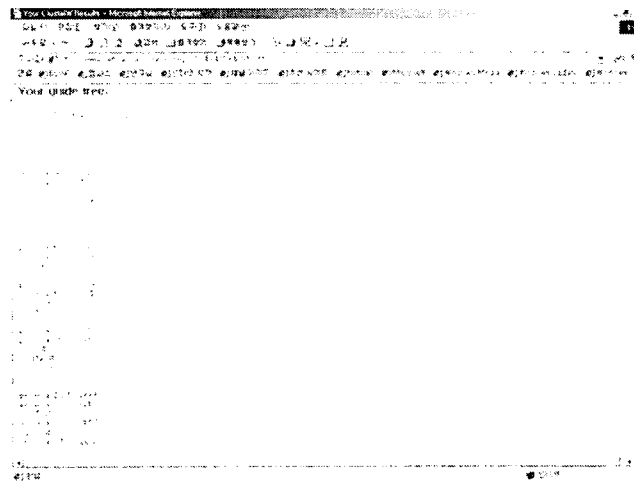
ClusterW (1) (www2.ebi.ac.uk/clusterw)

Input Fasta file



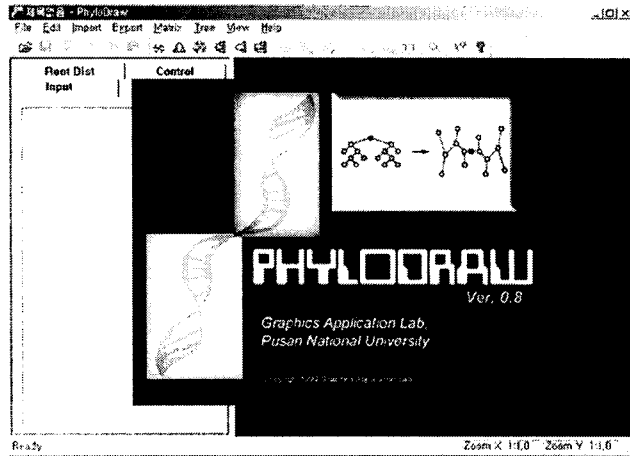
ClusterW(2) (www2.ebi.ac.uk/clusterw)

Output



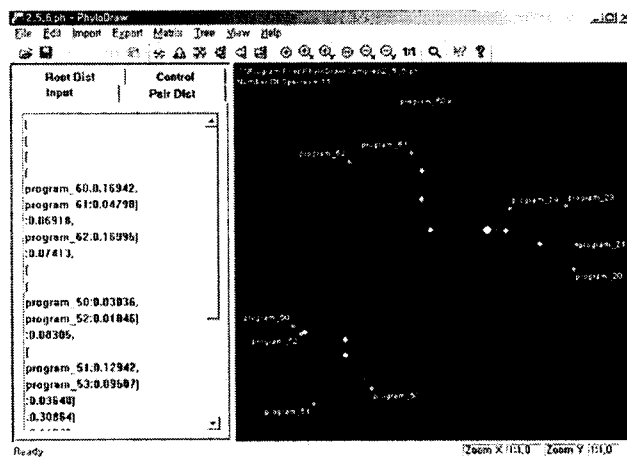
PhyloDraw

(cho et al, Bioinformatics 2001.)



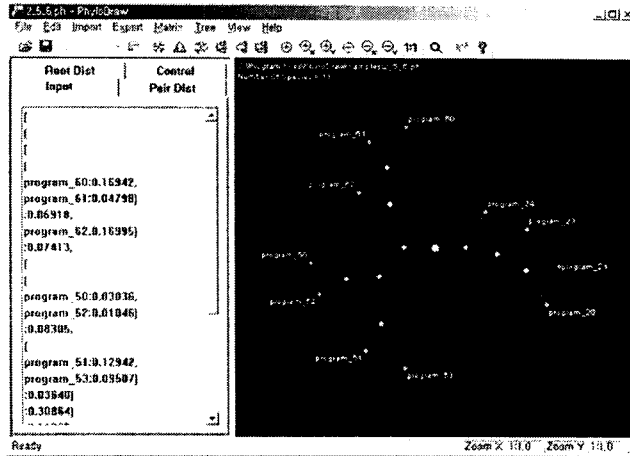
Experiment Result 1-0

□ 유사한 그룹을 이루는 11개의 프로그램



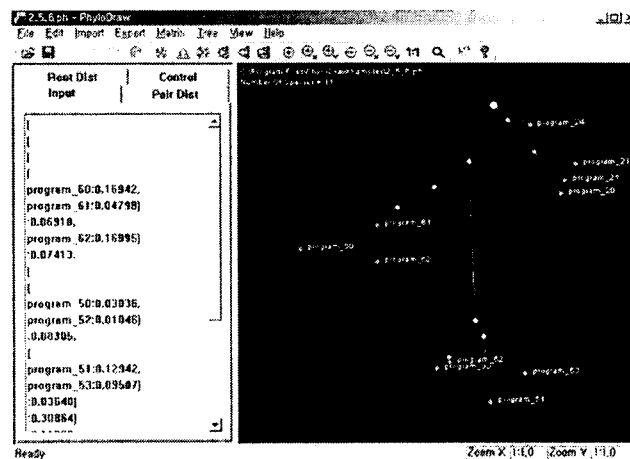
Experiment Result 1-1

□ 유사한 그룹을 이루는 **11**개의 프로그램



Experiment Result 1-2

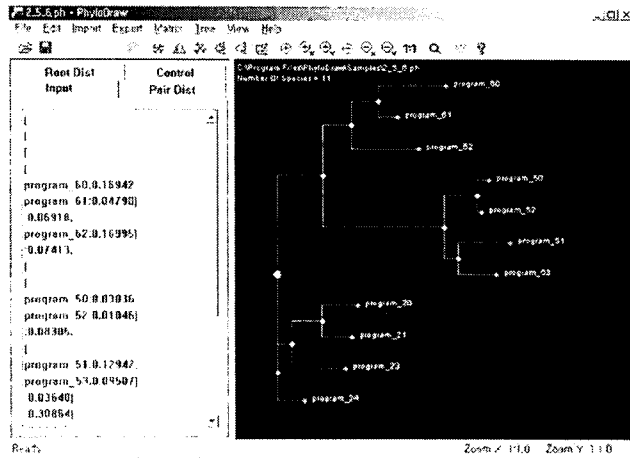
□ 유사한 그룹을 이루는 **11**개의 프로그램





Experiment Result 1-3

□ 유사한 그룹을 이루는 **11**개의 프로그램

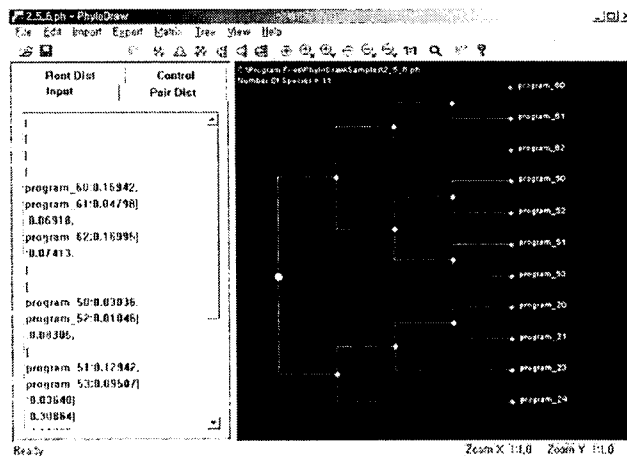


Graphics Application Lab



Experiment Result 1-4

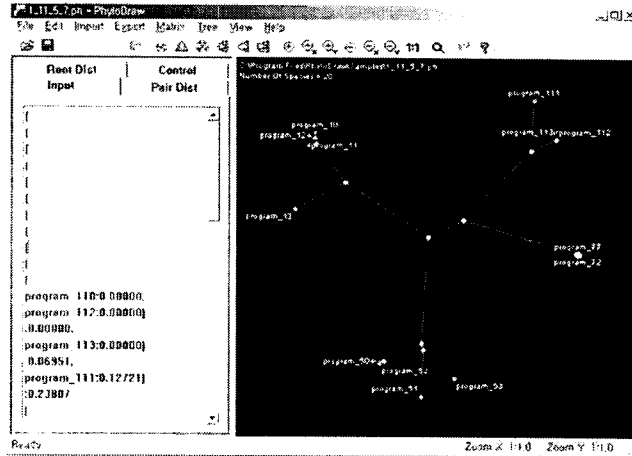
□ 유사한 그룹을 이루는 **11**개의 프로그램



Graphics Application Lab

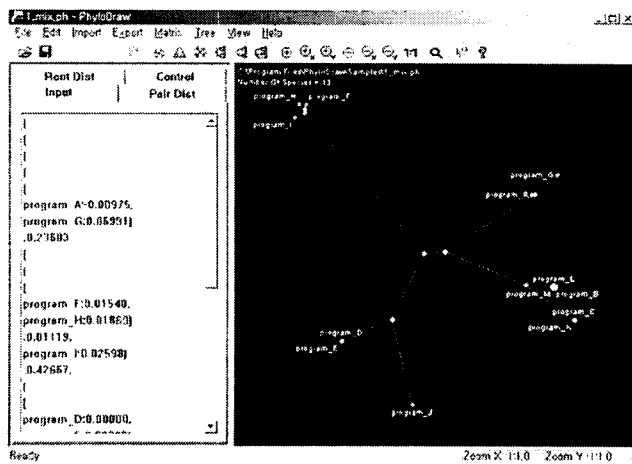
Experiment Result 2

□ 유사한 그룹을 이루는 13개의 프로그램



Experiment Result 3

□ 유사한 그룹을 이루는 13개의 프로그램



Experiment Result 4

□ 유사한 그룹을 이루는 **17**개의 프로그램

The screenshot shows the PhyloDraw interface. On the left, the 'Input' tab is active, displaying a list of programs and their distances from the root. On the right, a phylogenetic tree is displayed with 17 programs labeled as 'program_00' through 'program_16'. The tree shows a hierarchical structure where programs are grouped into clusters. The status bar at the bottom indicates 'Ready' and 'Zoom X: 1:1.0 Zoom Y: 1:1.0'.

Program	Distance
program_60	0.17175
program_61	0.04560
	0.06172
program_62	0.17741
	0.10489
program_10	0.01663
program_11	0.01746
	0.01126
program_12	0.02591
	0.09478

Experiment Result 5

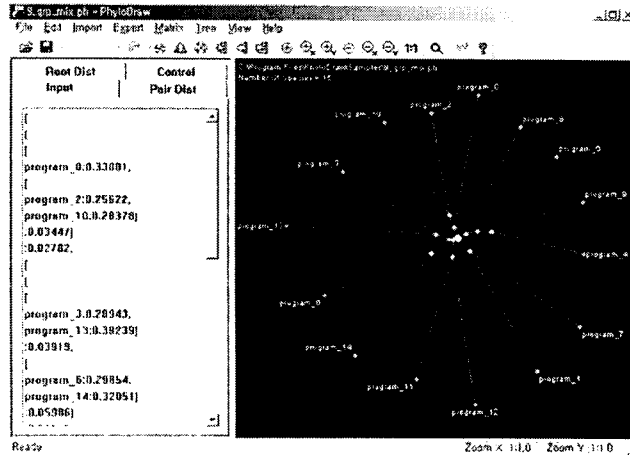
□ 유사한 그룹을 이루는 **21**개의 프로그램

The screenshot shows the PhyloDraw interface. On the left, the 'Input' tab is active, displaying a list of programs and their distances from the root. On the right, a phylogenetic tree is displayed with 21 programs labeled as 'program_00' through 'program_20'. The tree shows a hierarchical structure where programs are grouped into clusters. The status bar at the bottom indicates 'Ready' and 'Zoom X: 1:1.0 Zoom Y: 1:1.0'.

Program	Distance
program_60	0.16236
program_61	0.05593
	0.06135
program_62	0.17778
	0.14196
program_50	0.02564
program_52	0.01518
	0.09553
program_54	0.11935
program_53	0.10514
	0.02386
	0.24081

Experiment Result 6

- 유사도가 낮은 그룹을 이루는 **14**개의 프로그램



Conclusion

- **Alignment** 응용의 확장
- **Bioinformatics concept**의 활용
- **Linear Structure server** 구축
- **Document fingerprint**의 활용



Further Work

- ❑ **Program DNA-Bank server**
- ❑ **Copying Phylogenetics Building**
- ❑ **Multi-alignment**
- ❑ **Parametric Method**
 - o Fixed-size fingerprinting = program protein
 - o PAM for program copying behavior
 - o Real Practice
- ❑ **Korean-Report Oracle**
 - o Has this report some originality ? Then how much ?
- ❑ **Music Plagiarism** (patent, smallsoft, BIOventure)
 - o Melody, bit, harmonics sequence alignment



Research Problems

- ❑ **How to linearize a procedure call ?**
 - o Procedure call is a sort of directed graph
- ❑ **Block-based linear alignment**
 - o Tree alignment
- ❑ **Genomic sequence with fixed-size fingerprint**(may help a fast screening)



People in BIOINFORMATICS



Graphics Application Lab