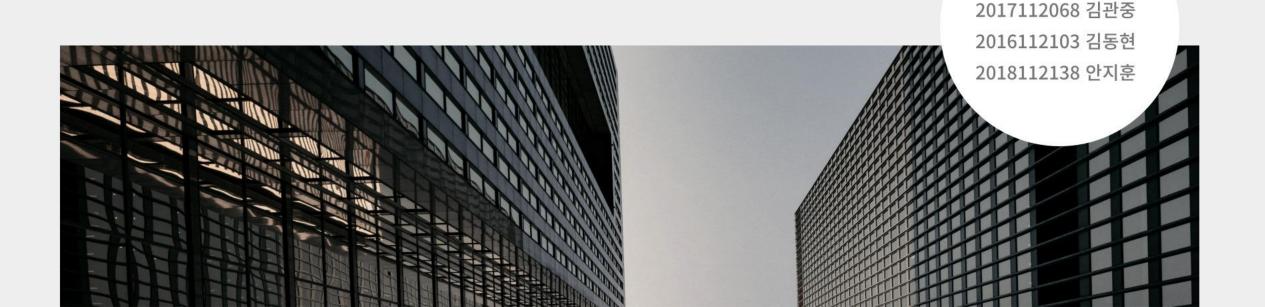


Team USB.

Multiple Sequence Alignment with Genetic Algorithm

With. Doungguk Univ.

2017112139 김규열





Contents

- 1. 주제 선정
- 2. Genetic Algorithm
- 3. SAGA (Sequecnce Alignment by Genetic Algorithm)
- 4. Experiment & Result



1. Global Alignment

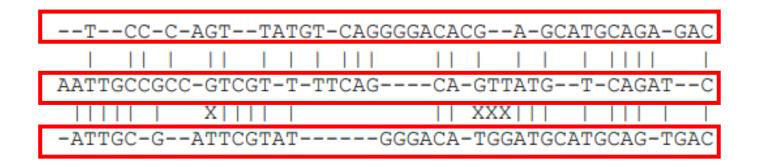
```
--T--CC-C-AGT--TATGT-CAGGGGACACG--A-GCATGCAGA-GAC
```

2. Local Alignment

```
tccCAGTTATGTCAGgggacacgagcatgcagagac
```

2개의 sequence에 대하여 alignment 수행 (Pairwise alignment) → 서열의 유사도 분석에 이용





3개 이상의 sequence에 대하여 alignment 수행



Pairwise Alignment보다 더 많은 생물학적 정보를 나타낼 수 있음



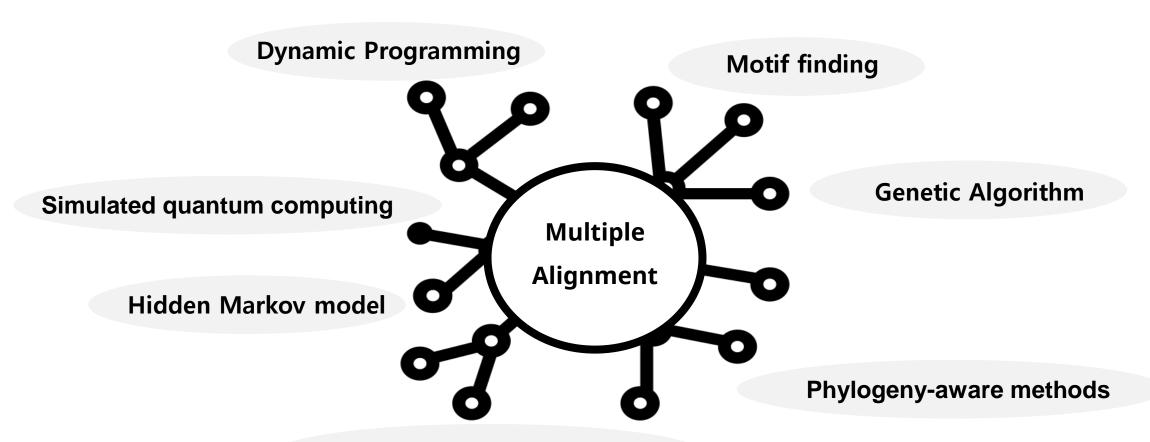
단백질, 서열 군에서 가변성 또는 보존 영역을 탐지



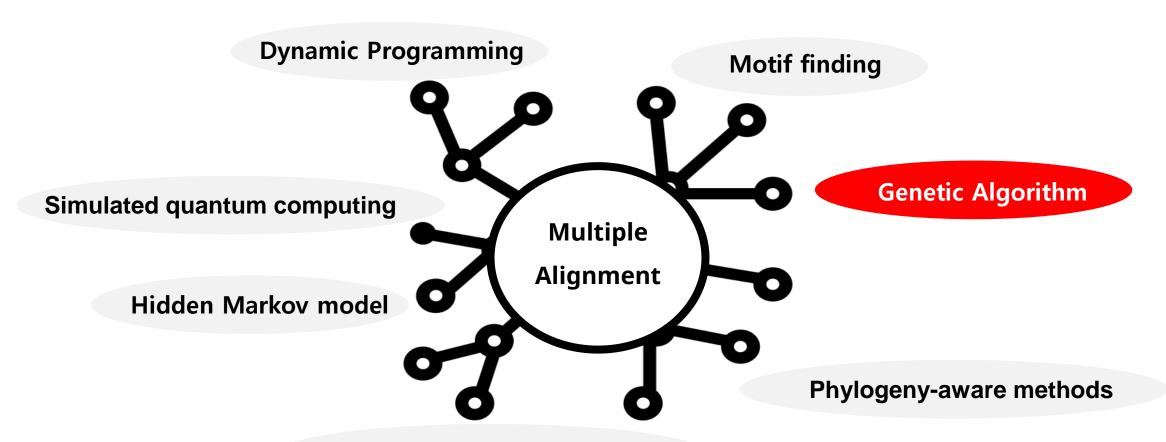
계통발생학적 분석을 통한 진화율 추정



기존 서열군 단백질 구조 예측 새롭게 서열화된 유전자 사이의 상동성 검출



Progressive alignment construction



Progressive alignment construction

Why? Genetic Algorithm..

1.

다윈의 진화론을 모방한 최적화 알고리즘

→ 바이오를 주제로 하는 수업 과 그 연관이 있다고 판단 2.

관련 논문, 저널 기사 존재

→ 프로젝트 구현, 결과 분석시 활용 가능 3.

다양한 분야, 특히 인공지능 분야에서 활발히 사용

→ 추후 다른 설계 프로젝트 에 활용 여지가 높음

Genetic Algorithm의 다양한 응용분야

응용 분야	응용 사례
최적화	수치적 함수의 최적화, 유전자 정보 해석 등
설계	비행기 날개의 공기 역학적 설계, 컴퓨터 통신망의 최적 설계 등
인공지능	문제해결 규칙의 자동 습득, 신경망 합성 및 학습, 패턴 인식, 자연언어 처리 등
시스템 분석 및 예측	시스템 동정, 케이오틱 시계열의 예측, 환율 변화 예측, 단백질 구조 분석 등
제어 및 로보틱스	이동 로봇의 경로 계획, 신경망 및 퍼지로직과 유전알고리즘의 결합에 의한 제어

Genetic Algorithm은 풀고자 하는 문제에 대한 가능해들을 일정한 형태의 자료구조(유전자)로 표현한 후, 점차 변형하여 더 나은 해를 만들어 내는 과정을 거침



"Multiple Sequence Alignment with Genetic Algorithm"

────〉 <Genetic Algorithm, MSA 연구 수행>

----> <Python을 활용한 개발 진행>

1. 진화 알고리즘을 사용한 복수 염기서열 정렬

The Korean Journal of Microbiology, Vol. 35, No. 2. June 1999, p. 115-120 Copyright: ⊙1997, The Microbiological Society of Korea

진화 알고리즘을 사용한 복수 염기서열 정렬

김 진^{1*} · 송민동² · 최흥식³ · 장연아³ 건국대학교 자연과학대학 전산과학과^{1*}, 분자생물학과², 한된대학교 컴퓨터공학부³

3개 이상의 DNA 혹은 단백질의 염기서열을 정렬하는 복수 염기서열 정렬(multiple sequence alignment)은 염기서열들 사이의 전화단계, gene regulation, 단백질의 구조와 기능에 관한 연구에 될수적인 도구이다. 복수 염기서열 점렬을 얻기 위한 기준의 방법은 progressive pairwise alignment와 잡이 빠른 실행시간 내에 만족할 만한 복수 염기서열 정렬을 제공하는 방법과, 최적의 복수 염기서열 정렬을 제공하는 실행시간이 상대적으로 21 dynamic programming과 같은 방법 등이 있다. 본 논문에 서는 진확 않고리즘을 사용하여 기준의 방법에서 제공하는 복수 염기서열 정렬을 짧은 시간 내에 보다 개선된 복수 염기서열 정렬을 활은 시간 내에 보다 개선된 복수 염기서열 의 자점을 보였다.

KEY WORDS ☐ multiple sequence alignment, genetic algorithm, dynamic programming, sequence comparison

생물학 역사상 가장 중요한 프로젝트의 하나인 Human Genome Project의 기본적인 목표는 인체의 게놈과 생명체의 유전자의 엄기서열인 인식을 목표로 하고 있다. 이 프로젝트에 의해 반생되는 엄청난 양의 염기서열 관련 테이터는 의약과 생물학 분야에 절대적인 역항력을 미치고 있으며 이러한 추세는 더욱 심화될 것이라 예상된다. 이러한 염기서열 관련 데이터를 처리하여 중요한 생물학적 정보를 얻기 위해서는 전산학의 도움이 필수적이다. 전산학에서 염기서열은 스트링으로 간주된다. 본 논문에서는 개놈 프로젝트에서 파생된 가장 중요한 문제 중에 하나인 복수 염기서열 정렬(multiple sequence alignment) 문제에 대하여 논한다(3-5, 7, 20).

여기시여 저편 © rbWzl(-----i-) TANA 미 DANA이 제무원

2. SAGA: Sequence Alignment by Genetic Algorithm

© 1996 Oxford University Press

Nucleic Acids Research, 1996, Vol. 24, No. 8 1515-1524

SAGA: sequence alignment by genetic algorithm

Cédric Notredame* and Desmond G. Higgins

EMBL outstation, The European Bioinformatics Institute, Hinxton Hall, Hinxton, Cambridge CB10 1RQ, UK

Received December 5, 1995; Revised and Accepted March 4, 1996

ABSTRACT

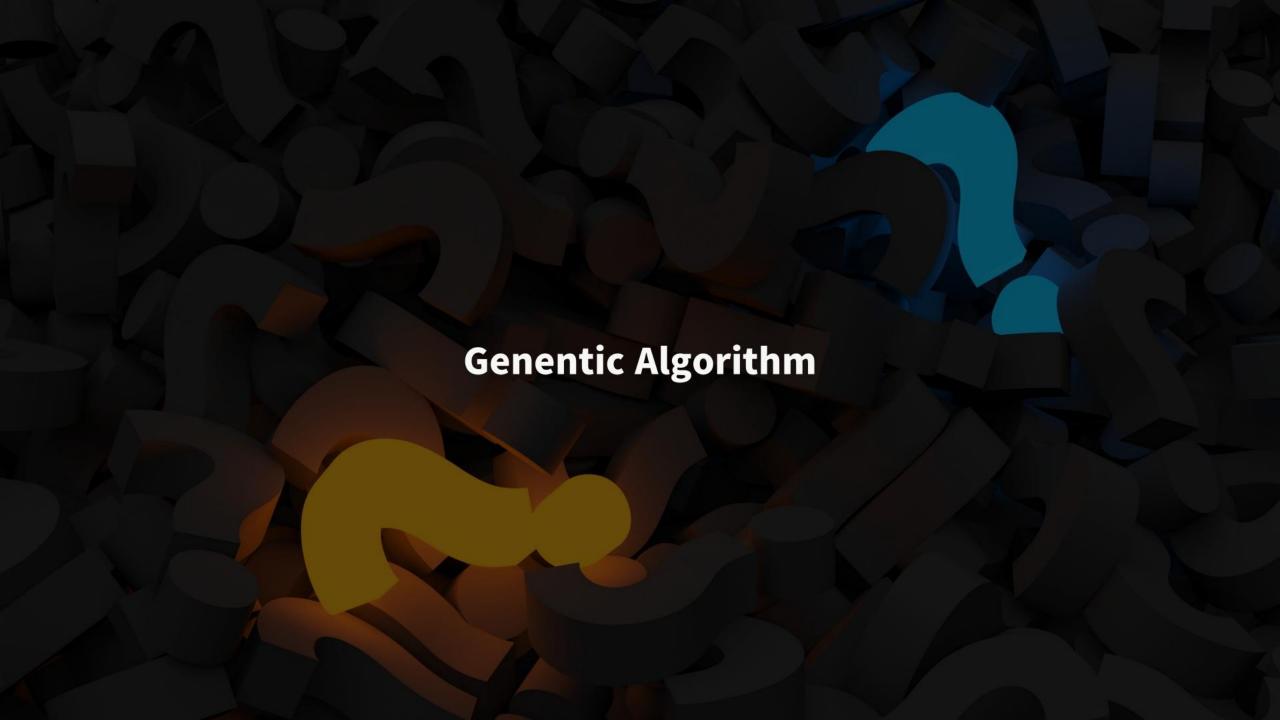
We describe a new approach to multiple sequence alignment using genetic algorithms and an associated software package called SAGA. The method involves evolving a population of alignments in a guasi evolutionary manner and gradually improving the fitness of the population as measured by an objective function which measures multiple alignment quality. SAGA uses an automatic scheduling scheme to control the usage of 22 different operators for combining alignments or mutating them between generations. When used to optimise the well known sums of pairs objective function, SAGA performs better than some of the widely used alternative packages. This is seen with respect to the ability to achieve an optimal solution and with regard to the accuracy of alignment by comparison with reference alignments based on sequences of known tertiary structure. The general attraction of the approach is the ability to optimise any objective function that one can invent.

INTRODUCTION

The simultaneous alignment of many nucleic acid or amino acid sequences is one of the most commonly used techniques in sequence analysis. Multiple alignments are used to help predict the secondary or tertiary structure of new sequences; to help demonstrate homology between new sequences and existing families; to help find diagnostic patterns for families; to suggest

There are two main alternatives to progressive alignment. One is to use hidden Markov models (HMMs; 5) which attempt to simultaneously find an alignment and a probability model of substitutions, insertions and deletions which is most self consistent. Currently, this approach is limited, in practice, to cases with very many sequences (e.g. 100 or more) but does have the great advantage of a sound link with probability analysis. A second approach is to use objective functions (OFs) which measure multiple alignment quality and to find the best scoring alignment. If the OF is well chosen or is an accurate measure of quality, then this approach has the advantage that one can be confident that the resulting alignment really is the best by some criterion. Unfortunately, the number of possible alignments which must be scored in order to choose the best one becomes astronomical for more than four or five sequences of reasonable length.

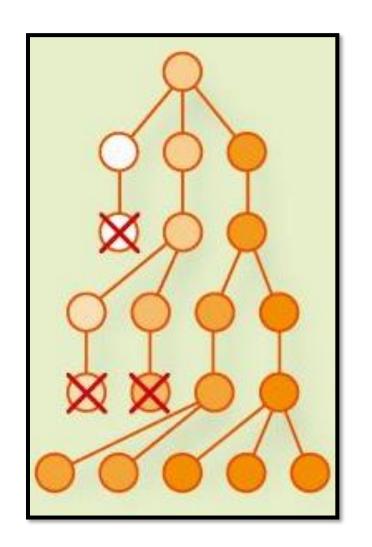
Two solutions to this problem exist. The MSA program (6,7) attempts to narrow down the solution space to a relatively small area where the best alignment is likely to be. It then guarantees finding the best alignment in this reduced space. Even with this reduction, it is limited to small examples of around seven or eight sequences at most. Nonetheless, it is the only method we know of that seems capable of finding the globally optimal alignment or close to it, starting with completely unaligned sequences. A second approach is to use stochastic optimisation methods such as simulated annealing (8), Gibbs sampling (9) or genetic algorithms (GAs; 10). Simulated annealing has been used on numerous occasions for multiple alignment (e.g. 11-13) but can be very slow and usually only works well as an alignment improver i.e. when the method is given an alignment that is already close to optimal and is not trapped in a local minimum. Gibbs sampling has been very successfully applied to the problem



!

Genetic Algorithm

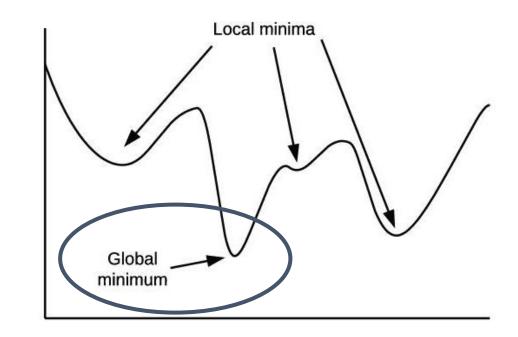
- 다윈의 <mark>자연선택 이론</mark>을 적용한 탐색, 최적화 알고리즘
- 문제에 대한 가능한 해들을 나열한 뒤, 교배, 변이를 통해 정확도가 높고 좋은 해를 생산하 게 된다
- 최적의 해를 찾는 방법론

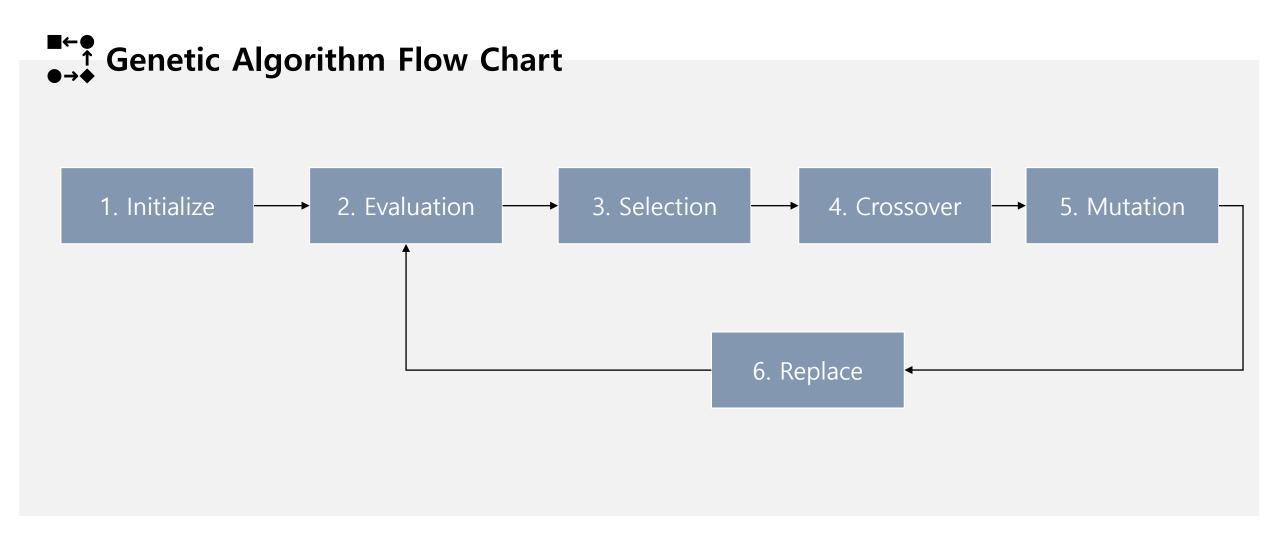


최적화 기법 중 전체 탐색영역에서 가장 좋은 해 를 찾는 것을 목적으로 하는 **전역 최적화 기법**

진화를 모방한 일종의 탐색 알고리즘

유전 알고리즘은 문제를 풀기위한 알고리즘이 기 보단 문제 풀이를 위한 **접근 방법**에 가까움





1. Initialize

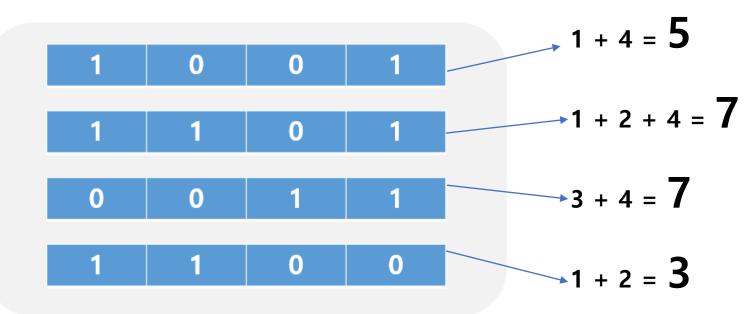
- 개체들로부터 모집단 생성 과정

1	0	0	1
1	1	0	1
0	0	1	1
1	1	0	0

2. Evaluation

- 모집단 내 각 개체에 대해 적합도 함수를 통해 적합도 측정
- → 일정 적합도 이상의 개체가 탄생하면, Termination

$$f(x) = \sum_{1}^{n} i \times arr[i]$$



적합도

- 평가하고자 하는 요소와 얼마나 부합하는지를 나타내는 척도

적합도 함수

- 현재 세대의 개체들을 평가하는 함수
- 다음 세대에 포함될 개체들을 선택하는 기준
- 유전 알고리즘을 적용하고자 하는 문제에 따라 다르게 설계

$$f(x) = \sum_{1}^{n} i \times arr[i]$$

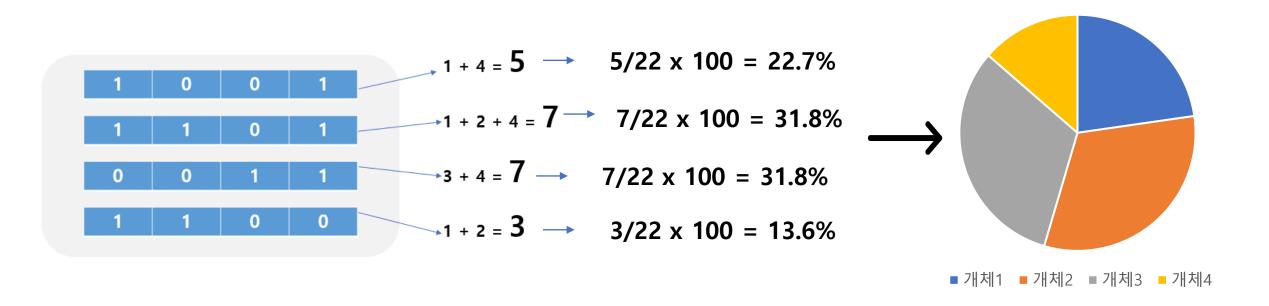
$$f(x) = \sum_{1}^{n} i \times arr[i]$$
 교체비용(A)= $\sum_{i=1}^{n-1} \sum_{j=1}^{n}$ 교체비용(S_i',S_j')

$$c(\alpha) = \sum_{0 < i < j \le n} c_2(\alpha_{ij}).$$

3. Selection

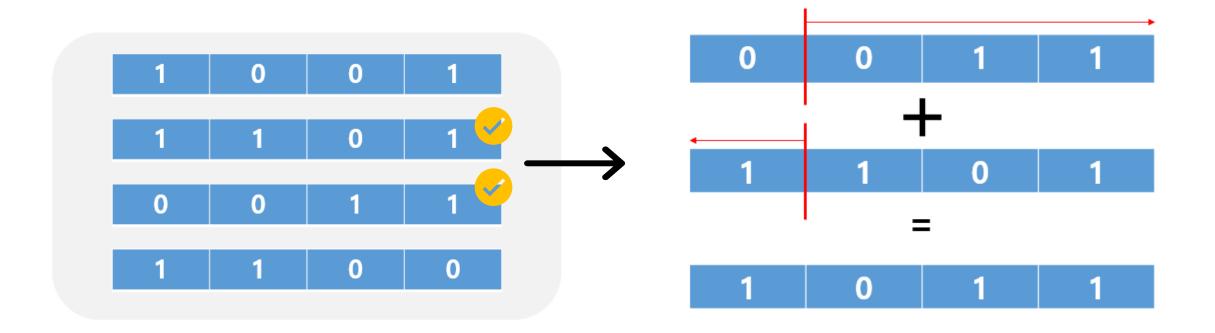
- 적합도에 따라 부여된 확률에 기반하여, <u>교배를 진행할 개체 선택</u> 하는 Roulette-wheel selection을 진행

Roulette

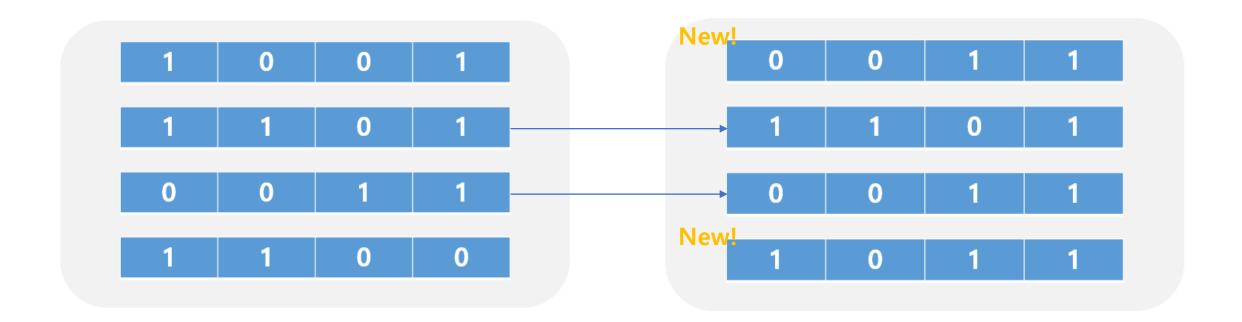


4. Crossover

- <u>선택된 두 개체의 일부를 교배(조합)</u>하여, 다음 세대의 모집단에 추가

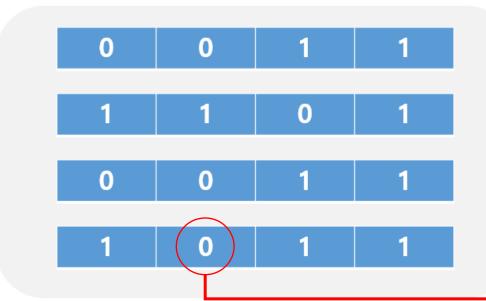


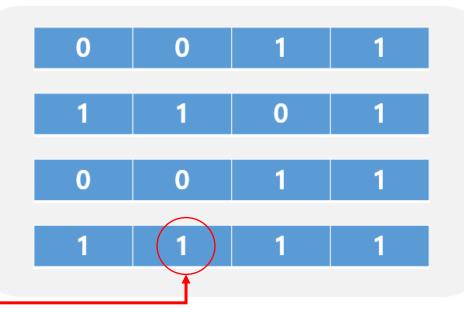
- 적합도 상위 50% 개체는 그대로 다음 세대로 전이, 하위 50% 개체는 교배를 통해 대체됨



5. Mutation

- 매우 낮은 확률로, 일정 개체에는 무작위 변이 진행

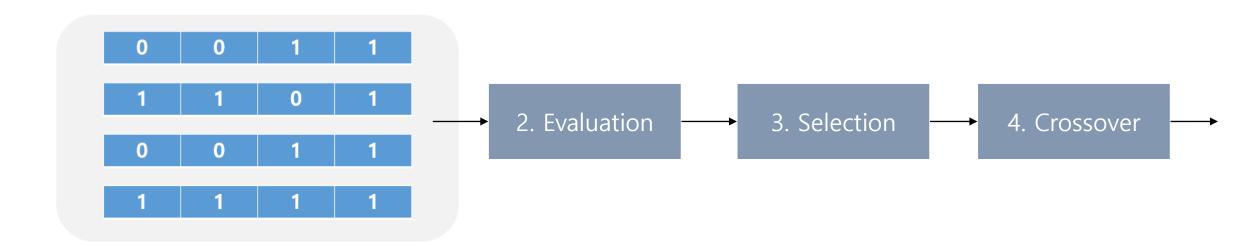


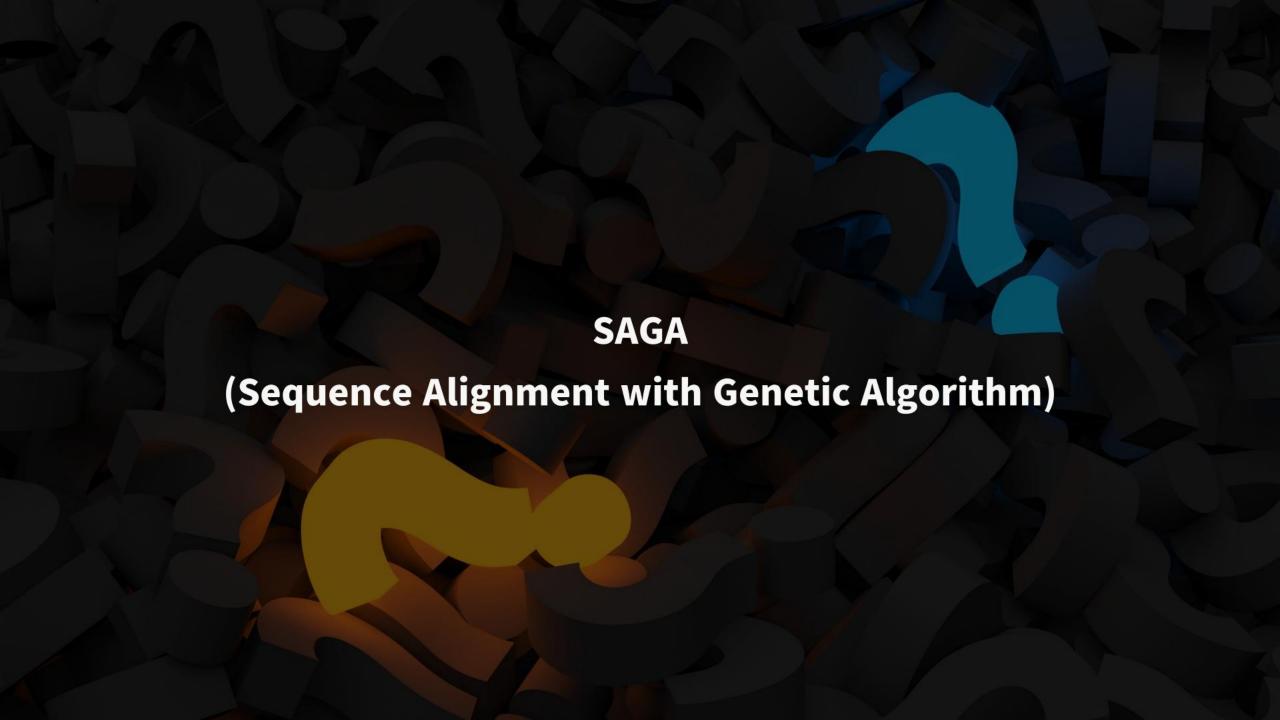


6. Replace

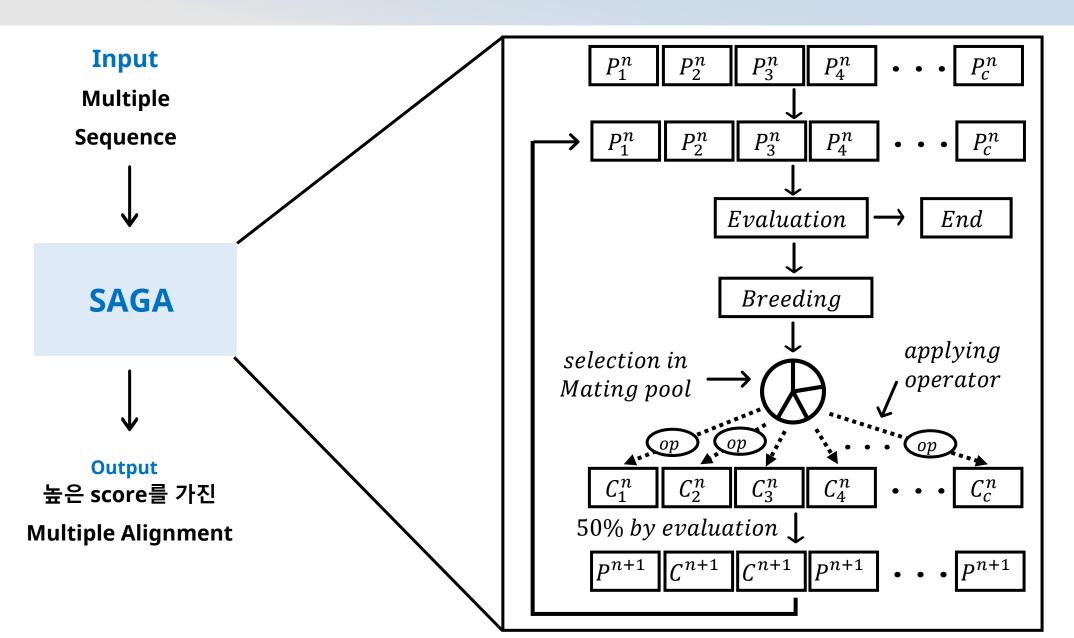
- 다음 세대로 전이

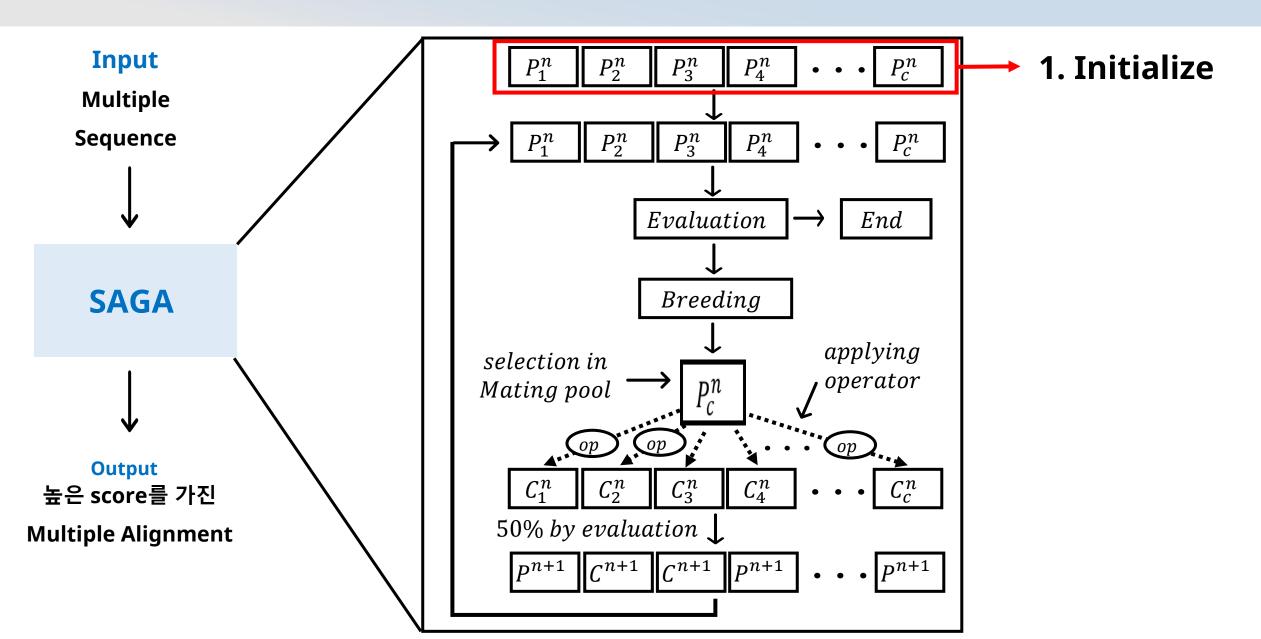
7. 2~6 단계 반복

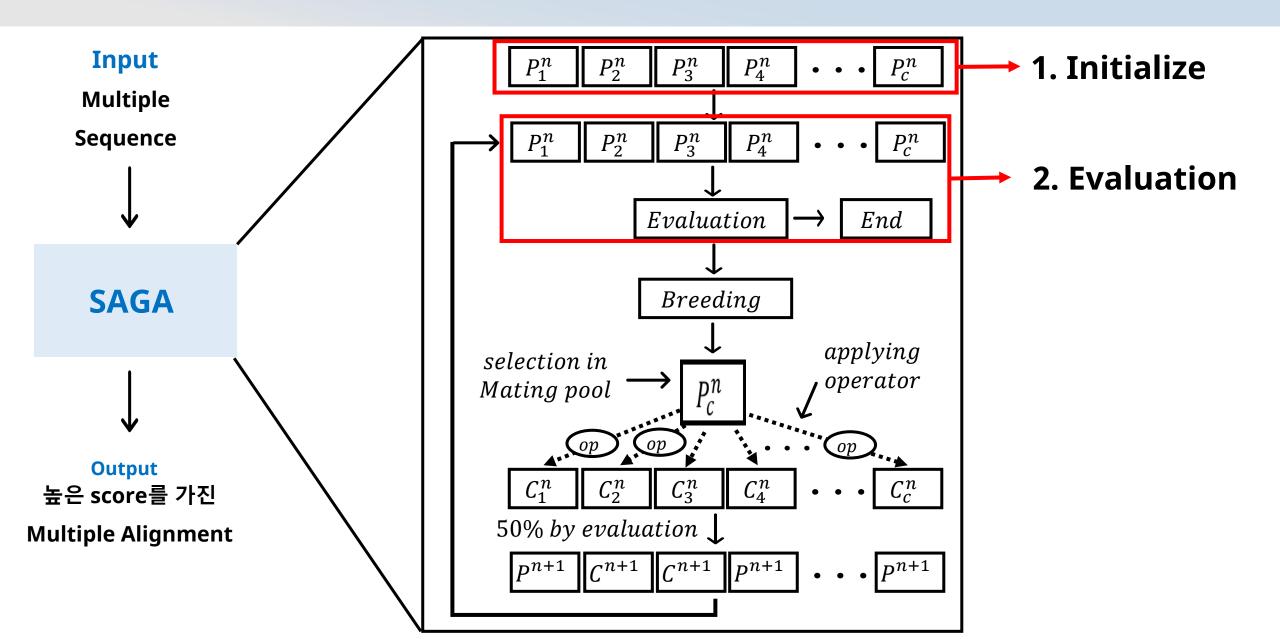


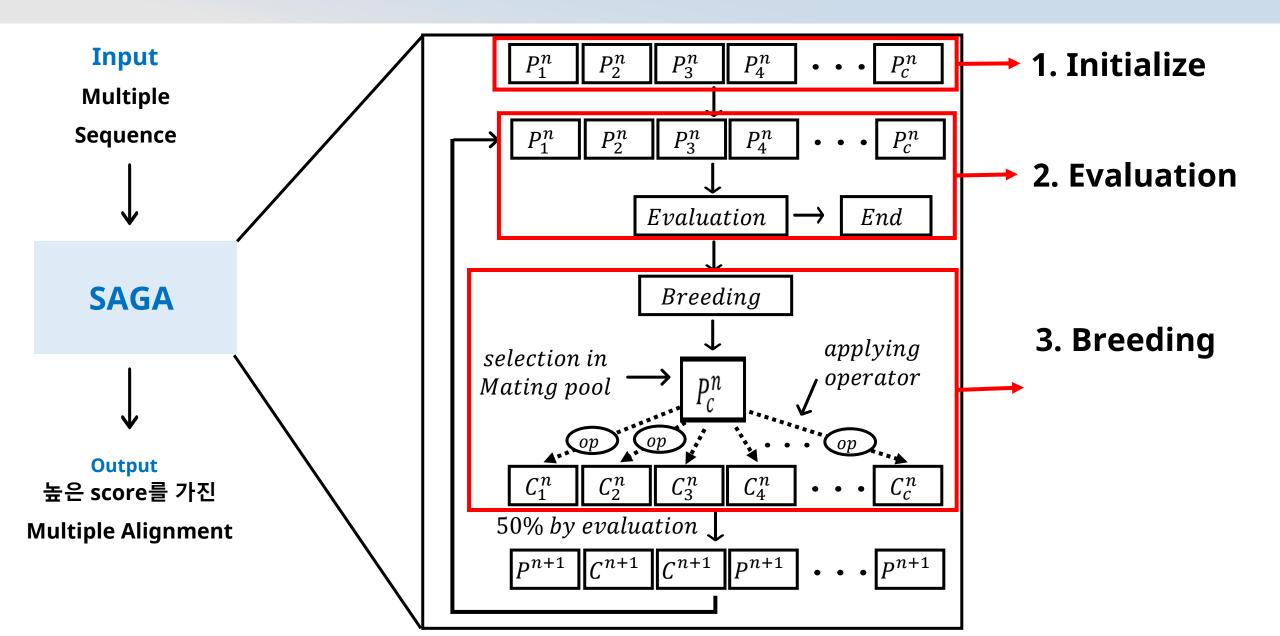


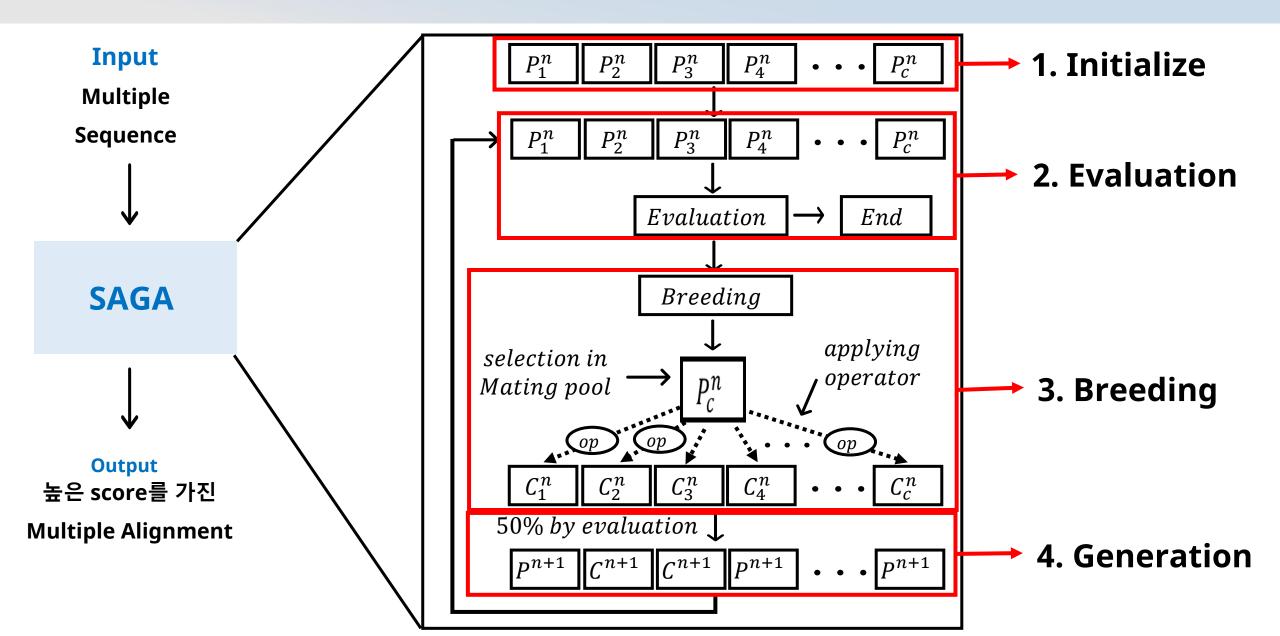


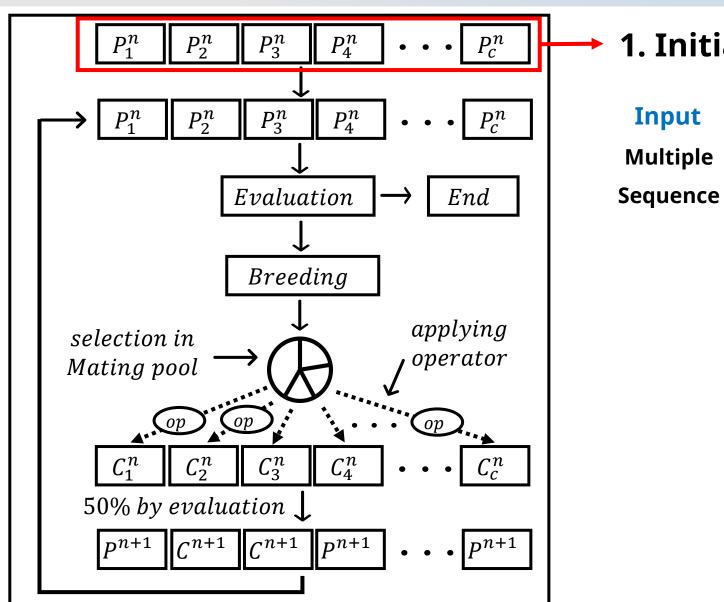






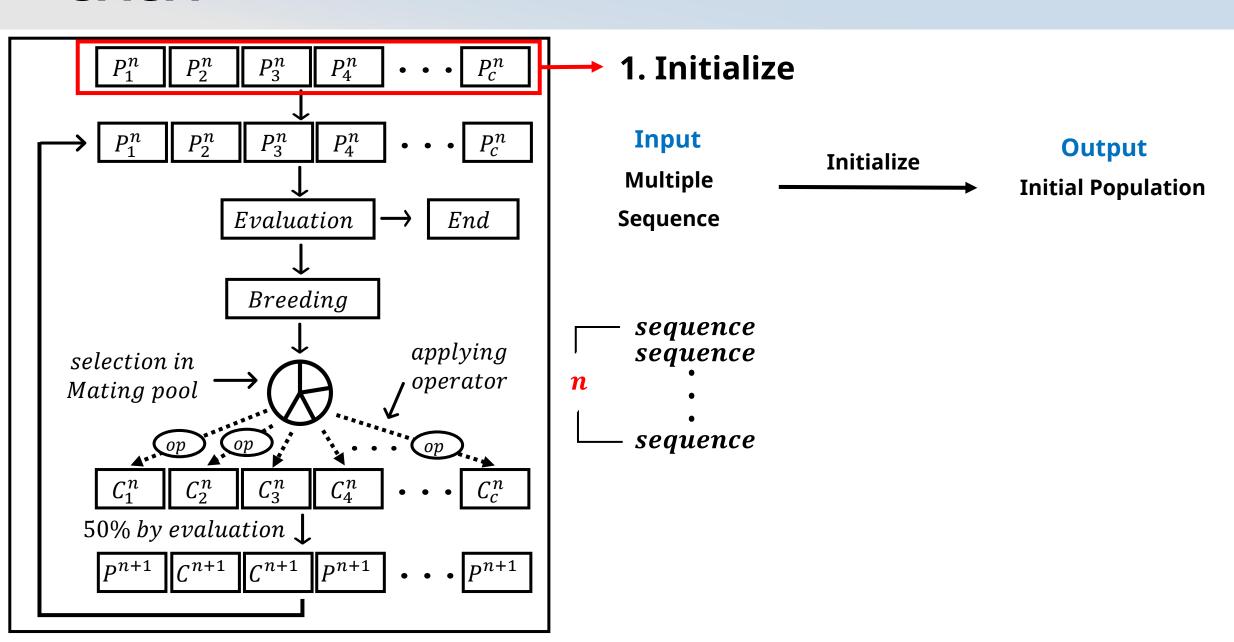


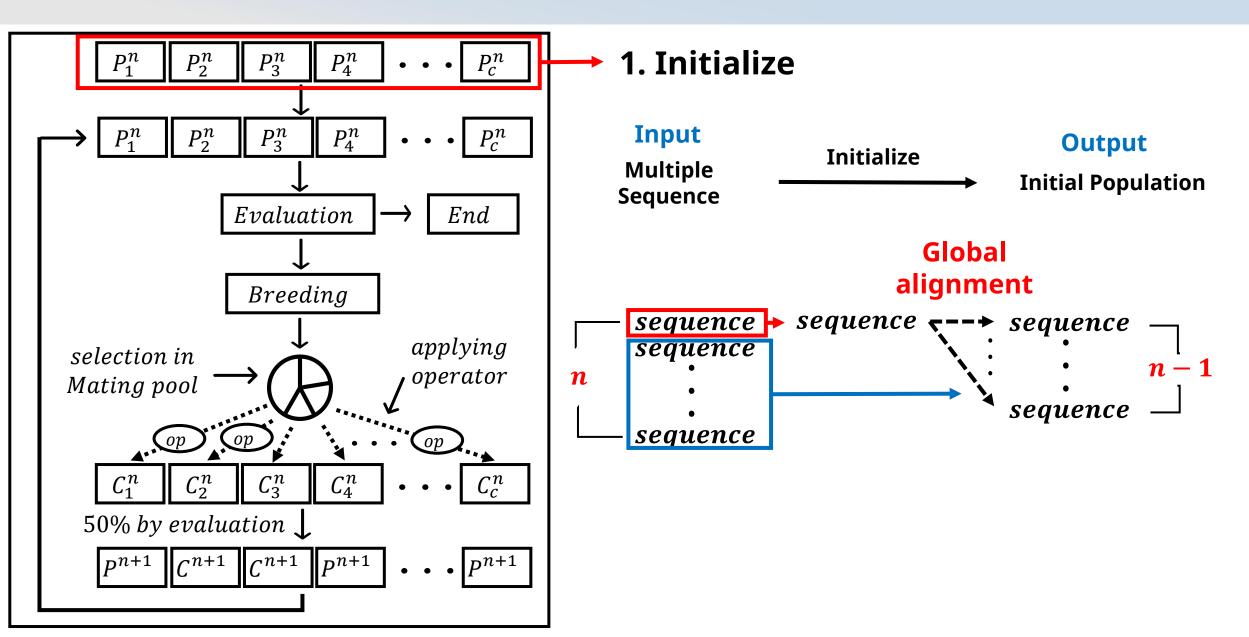


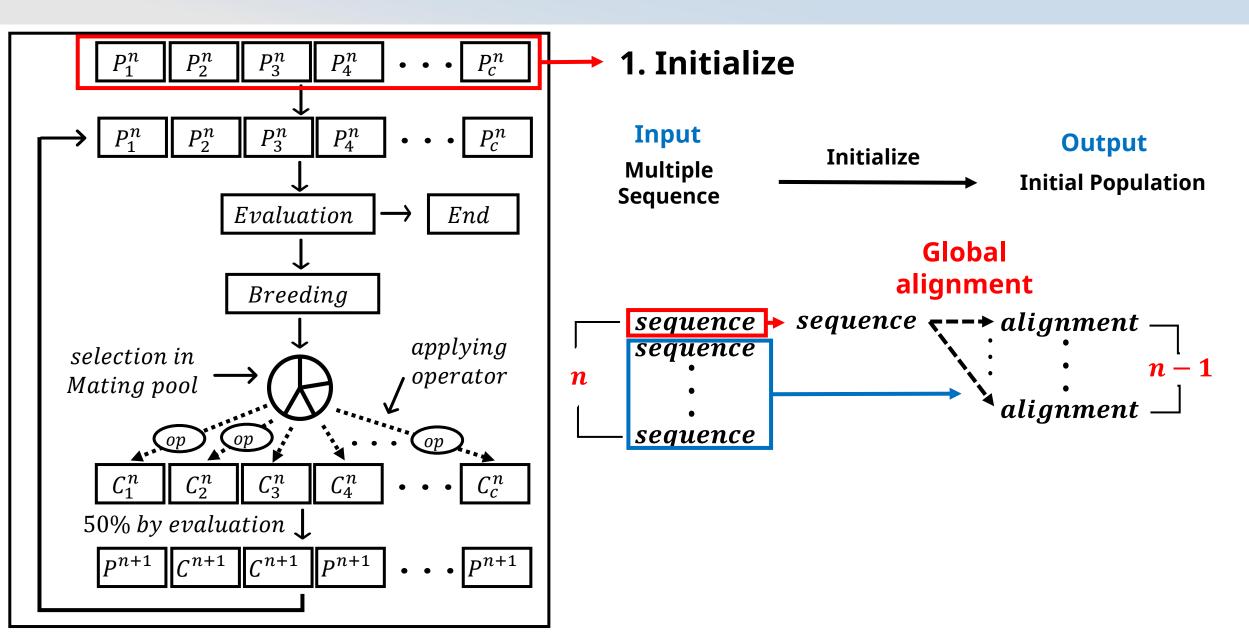


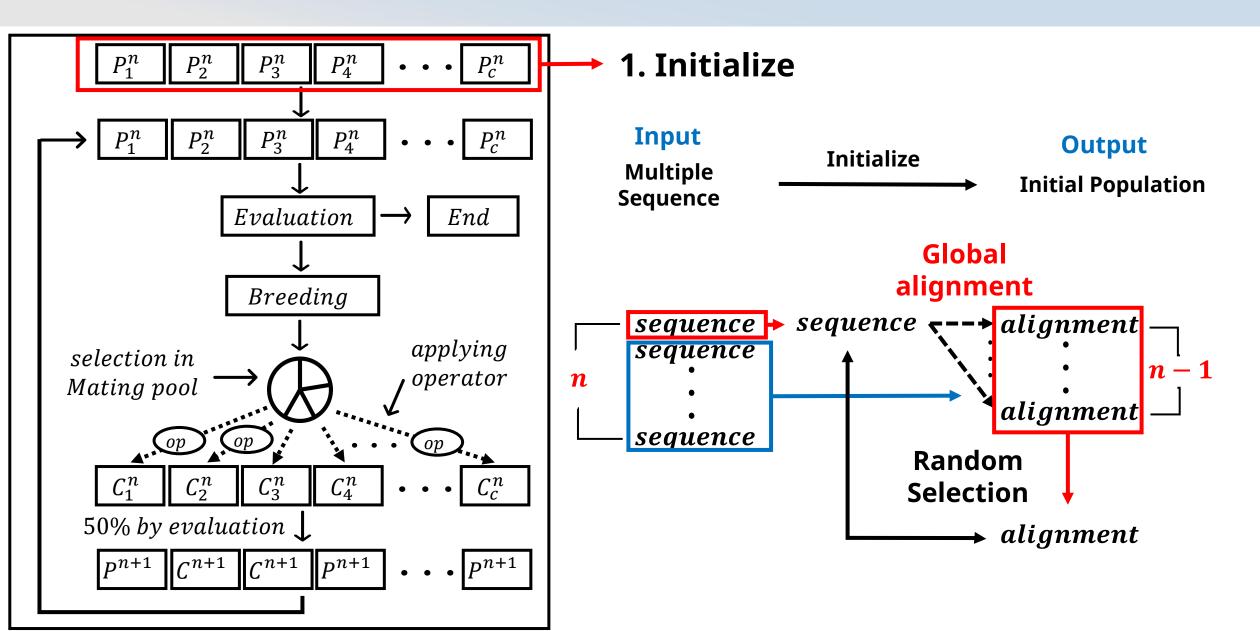
1. Initialize

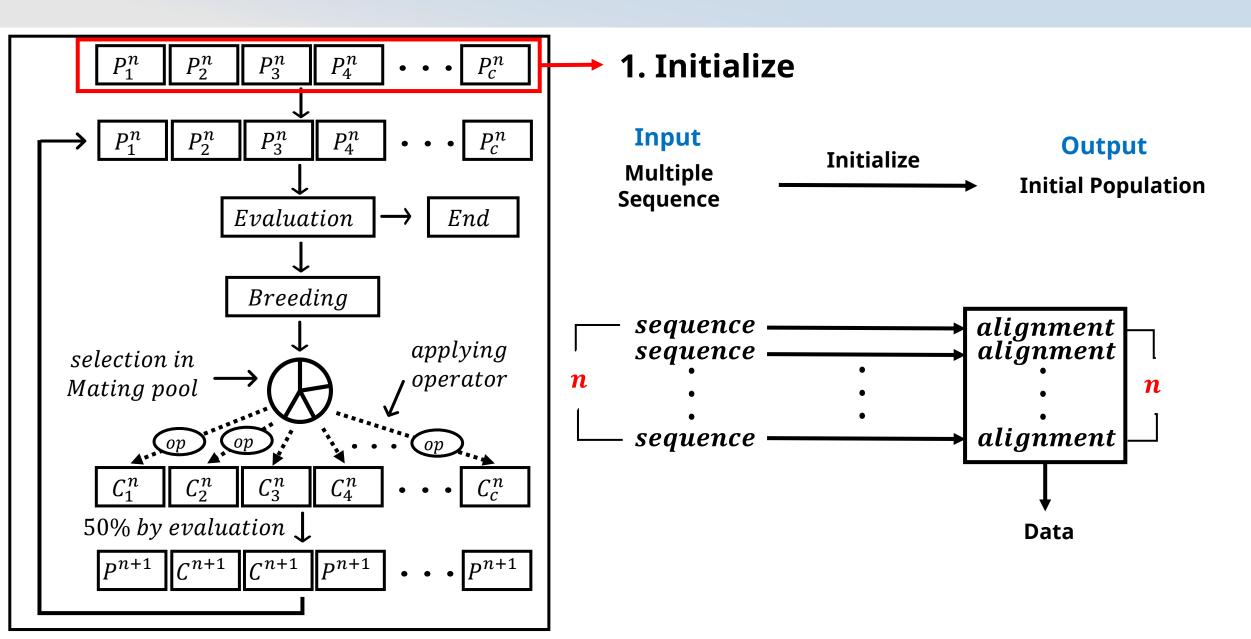
Output Initialize Initial Population

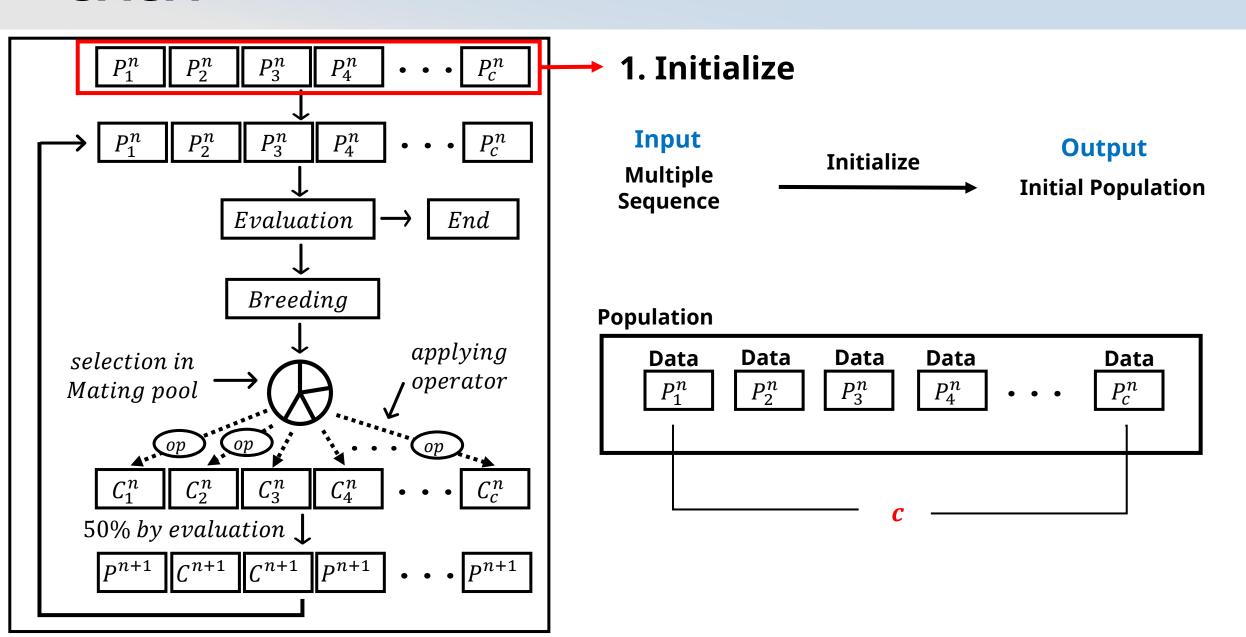


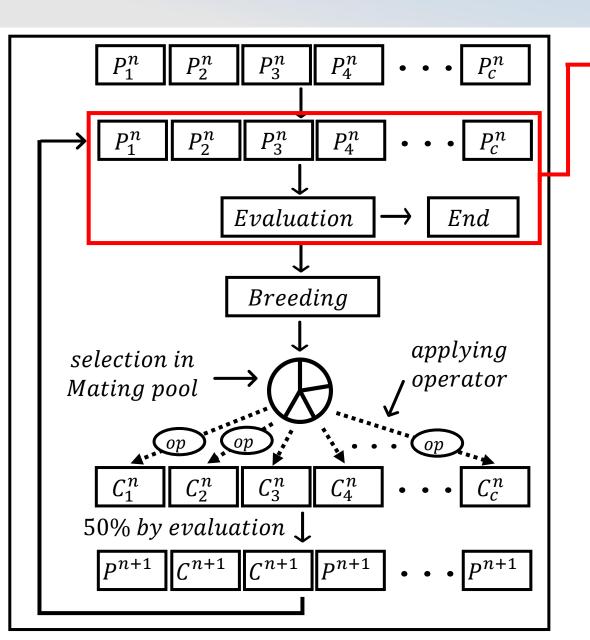












2. Evaluation

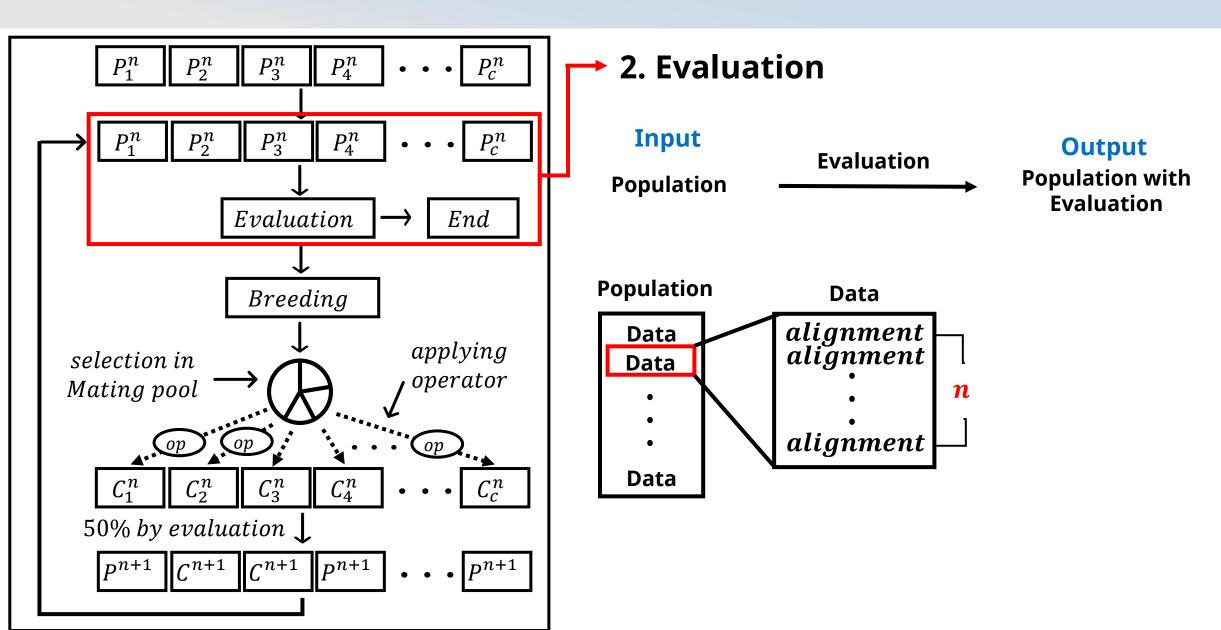
Input

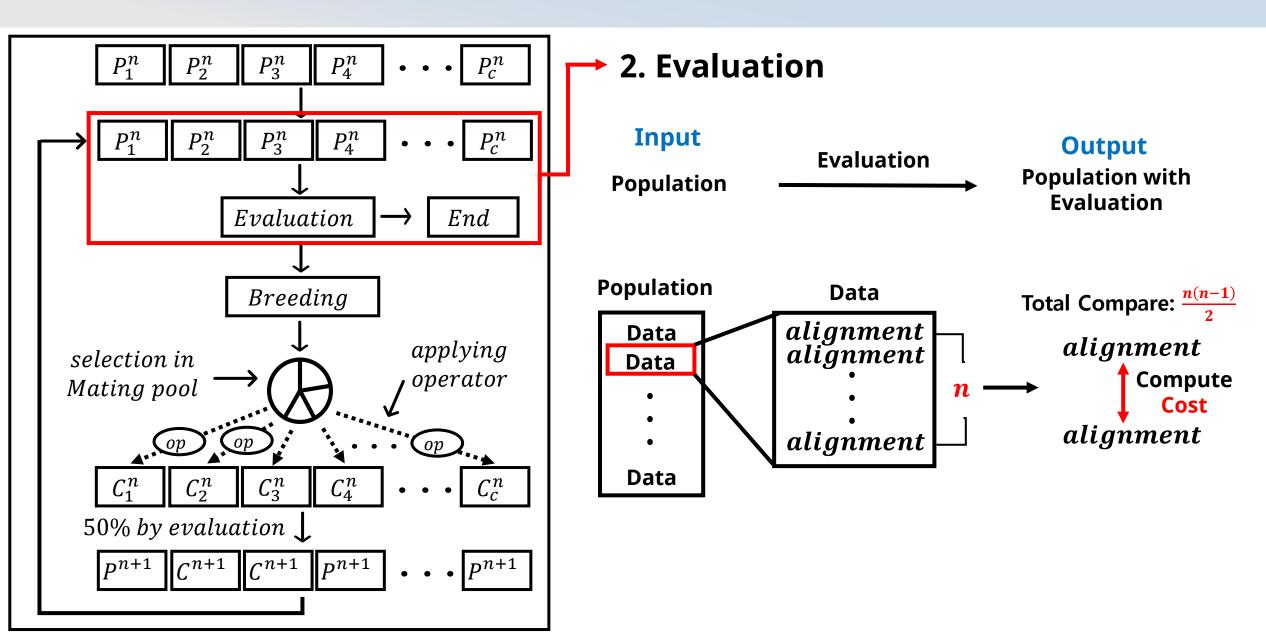
Evaluation

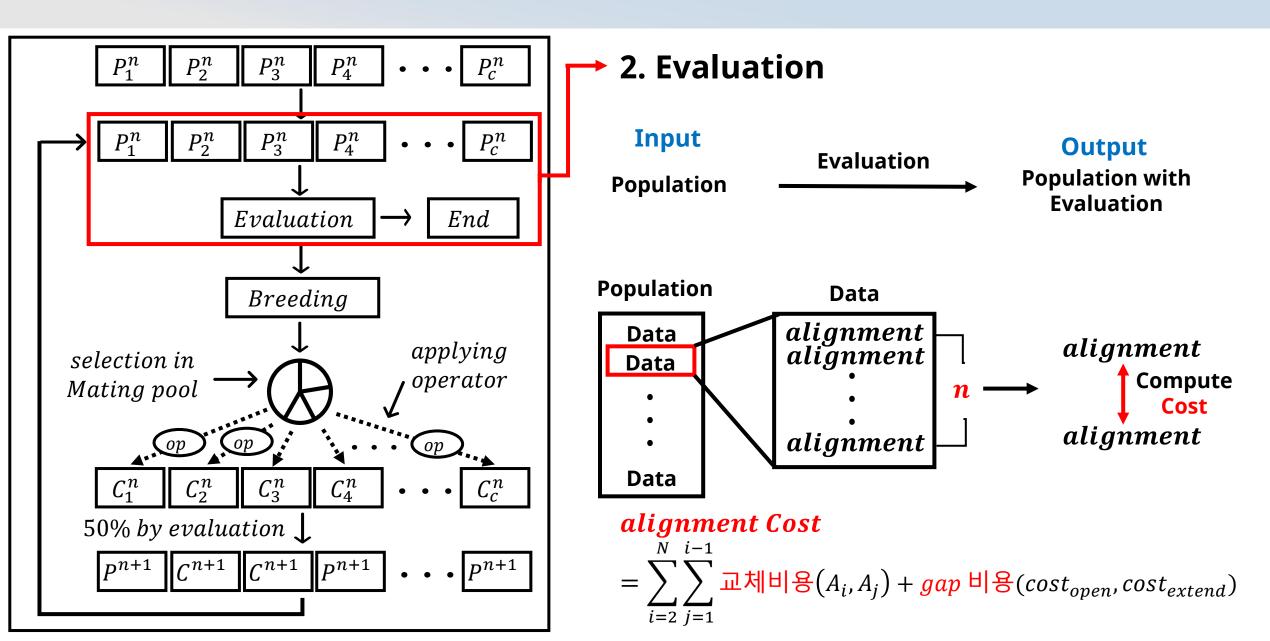
Population

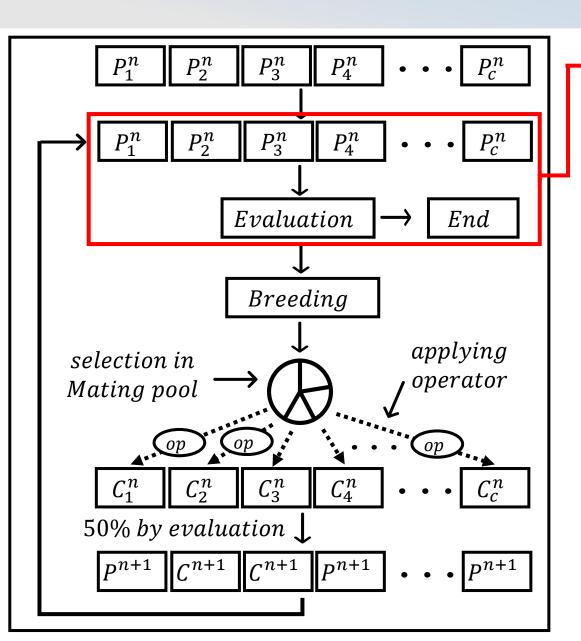
Evaluation

Evaluation









2. Evaluation

Input

Evaluation

Population

Evaluation

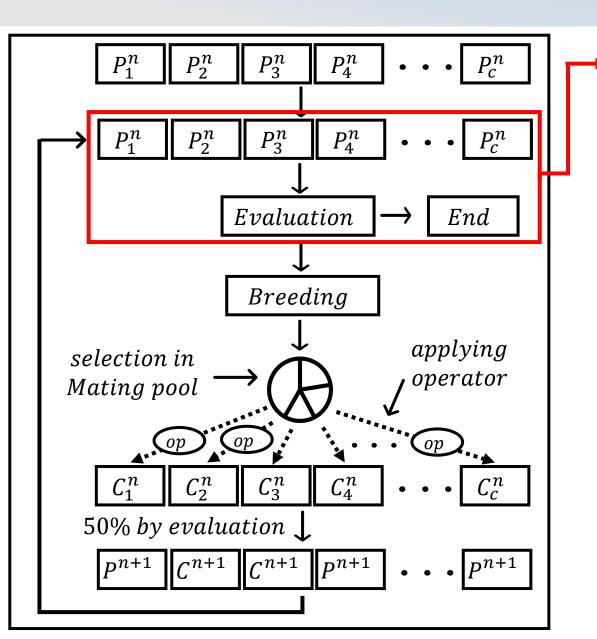
Evaluation

교체 비용

- alignment A가 alignment B로 바뀌기 위해 필요한 비용
- 이 비용은 BLOSUM 62테이블을 이용하여 계산 가능

alignment Cost

$$= \sum_{i=2}^{N} \sum_{j=1}^{i-1} \frac{1}{2} |A| + gap |A| + ga$$



2. Evaluation

Input

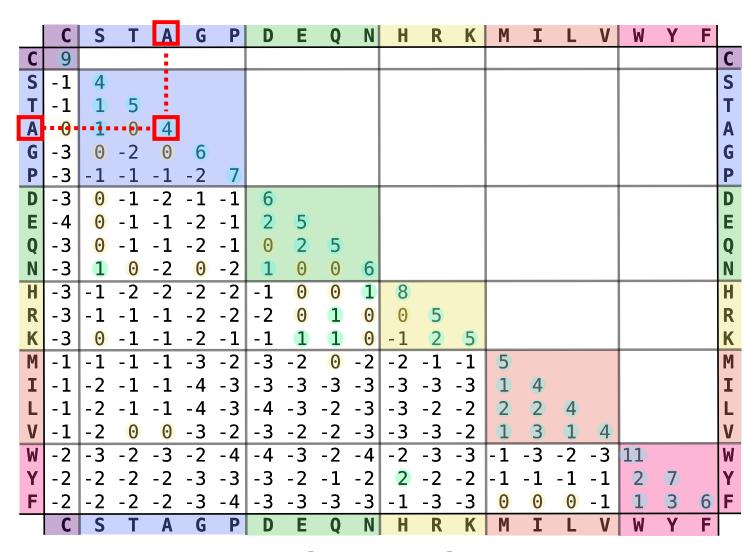
Ι	np	ut						Eva		at.	ior	•	Output											
op	si il	ati	on)					ııu	aι	101	1			Population with									
٥ŀ	, ai	uti	OI I												Evaluation									
	С	S	Т	Α	G	Р	D	Е	Q	N	Н	R	K	М	Ι	L	V	W	Υ	F				
C	9																				С			
	-1	4																			S			
Т	-1	1	5																		T			
Α	0	1	0	4																	Α			
G	-3	0	-2	0	6																G			
P	-3	-1	-1	-1	-2	7															P			
D	-3	0	-1	- 2	-1	-1	6														D			
Ε	-4	0	-1	-1	-2	-1	2	5													E			
Q	-3	0	-1	-1	-2	-1	0	2	5												Q			
N	-3	1	0	- 2	0	-2	1	0	0	6											N			
Н	-3	-1	-2	- 2	-2	-2	-1	0	0	1	8										Н			
R	-3	-1	- 1	- 1	-2	-2	-2	0	1	0	0	5									R			
K	-3	0	-1	- 1	-2	-1	-1	1	1	0	-1	2	5								K			
M	-1	-1	-1	- 1	-3	-2	-3	-2	0	-2	-2	- 1	- 1	5							M			
I	-1	-2	-1	- 1	-4	-3	-3	-3	-3	-3	-3	-3	-3	1	4						I			
L	-1	-2	- 1	- 1	-4	-3	-4	-3	-2	-3	-3	-2	- 2	2	2	4					L			
V	-1	-2	0	0	-3	-2	-3	-2	-2	-3		- 3	- 2	1	3	1	4				V			
W	-2	-3	_	- 3	-2	-4	-4	- 3	-2	-4	_	-3	-3	-1	-3	-2		11			W			
Y	-2	-2	-2	-2	-3	-3	-3	-2	-1	-2	2	-2	-2	-1	-1	- 1	-1	2	7		Y			
F	-2	-2	-2	-2		-4	-3	-3		-3	-1	-3	- 3	0	0	0	- 1	1	3	6	F			
	C	S	Т	Α	G	P	D	Е	Q	N	Н	R	K	M	Ι	L	V	W	Υ	F				
								[ΒL	.0	SU	М	62]										

1	С	S	Т	Α	G	Р	D	Е	Q	N	Н	R	K	M	I	L	V	W	Υ	F	
C	9																				C
S	- 1	4																			S
Т	- 1	1	5																		T
Α	0	1	0	4																	Α
G	-3	0	-2	0	6																G
Р	-3	-1	-1	-1	-2	7															P
D	-3	0	-1	-2	-1	-1	6														D
Ε	-4	0	- 1	- 1	- 2	-1	2	5													Ε
Q	-3	0	- 1	- 1	- 2	- 1	0	2	5												Q
N	-3	1	0	- 2	0	- 2	1	0	0	6											N
Н	-3	-1	- 2	- 2	- 2	- 2	-1	0	0	1	8										Н
R	-3	-1	- 1	- 1	-2	-2	-2	0	1	0	0	5									R
K	-3	0	- 1	- 1	-2	- 1	-1	1	1	0	- 1	2	5								K
M	-1	-1	- 1	-1	-3	-2	-3	-2	0	-2	-2	-1	-1	5							M
I	- 1	-2	- 1	- 1	-4	-3	-3	-3	-3	-3	-3	-3	-3	1	4						Ι
L	-1	-2	- 1	-1	-4	- 3	-4	-3	-2	-3	-3	-2	-2	2	2	4					L
V	- 1	-2	0	0	-3	-2	-3	-2	-2	-3	-3	-3	-2	1	3	1	4				٧
W	-2	-3	-2	-3	-2	- 4	- 4	-3	-2	-4	-2	-3	-3	-1	-3	-2	- 3	11			W
Υ	-2	-2	-2	-2	-3	- 3	-3	-2	- 1	-2	2	-2	-2	- 1	-1	-1	- 1	2	7		Υ
F	-2	- 2	- 2	-2	-3	- 4	-3	-3	-3	-3	-1	-3	-3	0	0	0	- 1	1	3	6	F
	C	S	Т	Α	G	Р	D	Е	Q	N	Н	R	K	M	Ι	L	V	W	Υ	F	

Alignment

1: ATGCGGTT

2: AGTCACGT



Alignment

	С	S	Т	Α	G	Р	D	Е	Q	N	Н	R	K	M	Ι	L	V	W	Υ	F	
C	9																				C
S	-1	4																			S
Т	- 1	1	5																		Т
Α	0	1	0	4																	Α
G	· - 3 ·	0-	-2	0	6																G
Р	-3	-1	- 1	-1	-2	7															Р
D	-3	0	- 1	-2	-1	- 1	6														D
Е	-4	0	- 1	- 1	-2	- 1	2	5													Ε
Q	-3	0	- 1	- 1	- 2	- 1	0	2	5												Q
N	-3	1	0	-2	0	- 2	1	0	0	6											N
Н	-3	-1	- 2	-2	- 2	- 2	-1	0	0	1	8										Н
R	-3	-1	- 1	- 1	- 2	- 2	-2	0	1	0	0	5									R
K	-3	0	- 1	- 1	-2	- 1	- 1	1	1	0	-1	2	5								K
M	- 1	- 1	-1	-1	-3	-2	-3	- 2	0	-2	- 2	-1	-1	5							М
I	- 1	- 2	- 1	-1	- 4	- 3	-3	- 3	-3	-3	-3	-3	-3	1	4						Ι
L	- 1	- 2	- 1	- 1	-4	-3	-4	-3	-2	-3	- 3	-2	-2	2	2	4					L
V	- 1	- 2	0	0	-3	-2	-3	-2	-2	-3	-3	-3	- 2	1	3	1	4				٧
W	-2	-3	- 2	-3	-2	- 4	- 4	- 3	-2	-4	-2	-3	-3	-1	-3	-2	- 3	11			W
Υ	-2	-2	- 2	-2	-3	- 3	-3	- 2	-1	-2	2	-2	-2	- 1	-1	- 1	- 1	2	7		Υ
F	-2	-2	-2	-2	-3	-4	-3	-3	-3	-3	- 1	-3	-3	0	0	0	-1	1	3	6	F
	С	S	Т	Α	G	P	D	Е	Q	N	Н	R	K	M	Ι	L	V	W	Υ	F	

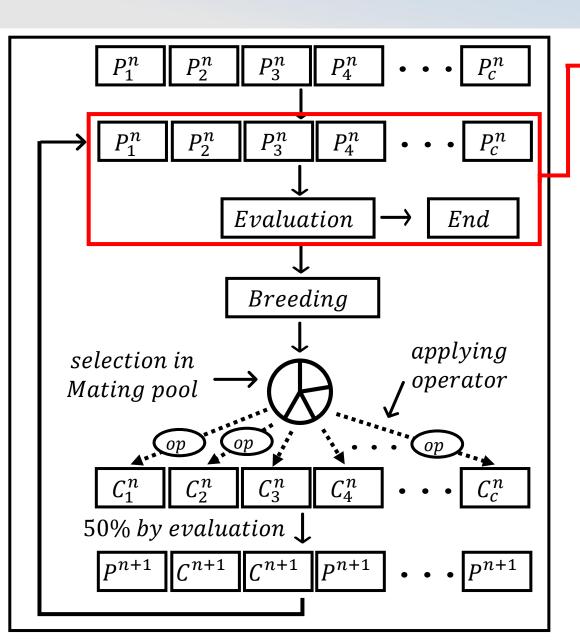
Alignment

[С	S	Т	Α	G	Р	D	Е	Q	N	Н	R	K	М	Ι	L	V	W	Υ	F	
С	9																				C
S	-1	4																			S
Т	-1	1	5																		T
Α	0	1	0	4																	Α
G	-3	0	-2	0	6																G
P	-3	-1	-1	-1	-2	7															P
D	-3	0	- 1	- 2	- 1	- 1	6														D
Е	-4	0	- 1	- 1	-2	- 1	2	5													Ε
Q	-3	0	- 1	-1	-2	- 1	0	2	5												Q
N	-3	1	0	-2	0	-2	1	0	0	6											N
Н	-3	-1	- 2	-2	-2	-2	-1	0	0	1	8										Н
R	-3	-1	- 1	- 1	-2	-2	- 2	0	1	0	0	5									R
K	-3	0	- 1	- 1	- 2	- 1	- 1	1	1	0	-1	2	5								K
M	- 1	- 1	- 1	- 1	- 3	- 2	-3	- 2	0	-2	-2	- 1	- 1	5							M
I	- 1	-2	- 1	- 1	- 4	- 3	-3	-3	-3	-3	- 3	-3	-3	1	4						Ι
L	-1	-2	- 1	-1	- 4	- 3	-4	-3	-2	-3	- 3	-2	-2	2	2	4					L
V	-1	-2	0	0	-3	- 2	-3	- 2	-2	-3	- 3	-3	-2	1	3	1	4				٧
W	-2	-3	-2	-3	-2	- 4	-4	-3	-2	-4	-2	-3	-3	- 1	-3	-2	- 3	11			W
Υ	-2	-2	-2	- 2	- 3	- 3	-3	-2	-1	-2	2	-2	-2	- 1	-1	- 1	-1	2	7		Y
F	-2	- 2	-2	-2	- 3	- 4	-3	-3	-3	-3	- 1	-3	-3	0	0	0	-1	1	3	6	F
	C	S	Т	Α	G	P	D	Е	Q	N	Н	R	K	M	Ι	L	V	W	Υ	F	

Alignment

1: ATGCGGTT

2: A G T C A C G T



2. Evaluation

Input

Evaluation

Population With

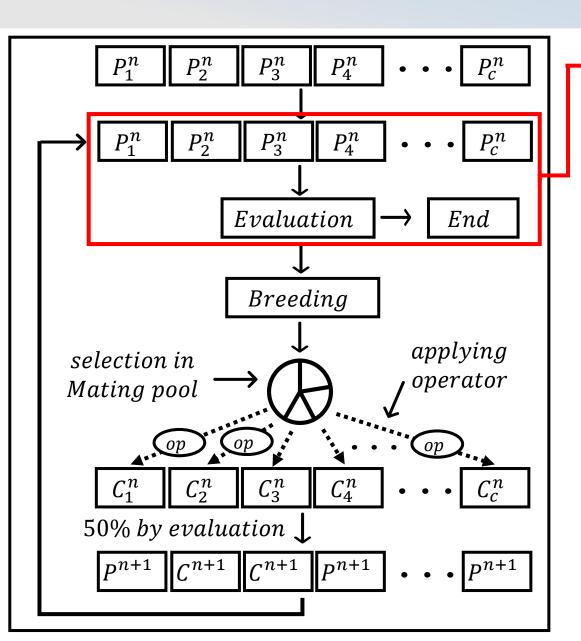
Evaluation

Gap 비용

- Gap 비용은 alignment의 기호 중 '一'에 부과되는 비용이다.
- 여기서 사용된 gap 비용은 natural affine gap 비용이다.

alignment Cost

$$= \sum_{i=2}^{N} \sum_{j=1}^{i-1} \frac{1}{2} \frac{1}{2}$$



2. Evaluation

Input

Evaluation

Population With

Evaluation

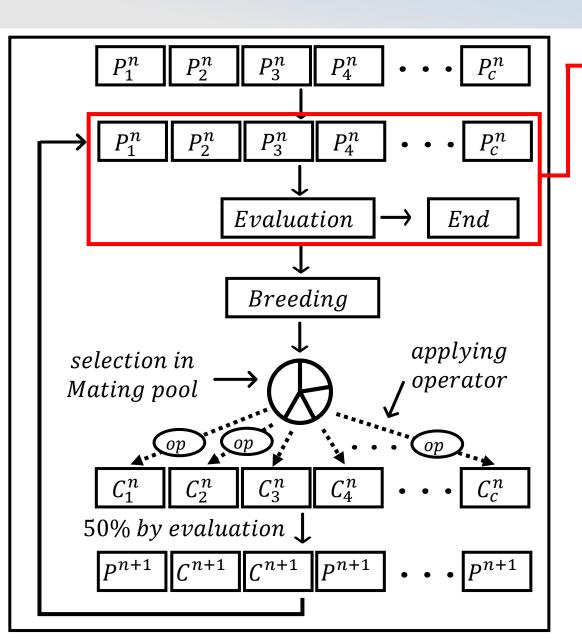
Gap 비용

- Gap 비용은 alignment의 기호 중 '一'에 부과되는 비용이다.
- 여기서 사용된 gap 비용은 natural affine gap 비용이다.

 $cost_{open}$: gap이 처음 나왔을 때 주는 비용 $cost_{extend}$: gap이 나올 때 마다 주는 비용

alignment Cost

$$= \sum_{i=2}^{N} \sum_{j=1}^{i-1} \frac{1}{2} \frac{1}{2}$$



2. Evaluation

Population Evaluation Population Evaluation

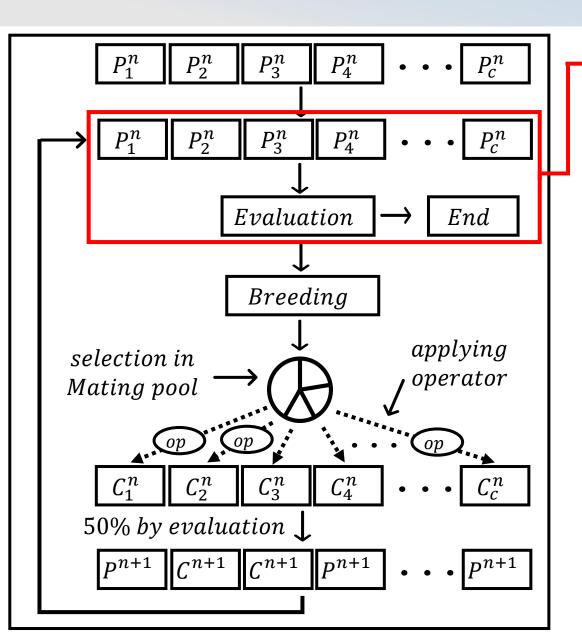
Gap 비용

- Gap 비용은 alignment의 기호 중 '一'에 부과되는 비용이다.
- 여기서 사용된 gap 비용은 natural affine gap 비용이다.

cost_{open}: -1 cost_{extend}: -0.5라 가정

4-1.5

$$cost_{open} + cost_{extend} = -1 - 0.5 = -1.5$$



2. Evaluation

Population Evaluation Population Evaluation

Gap 비용

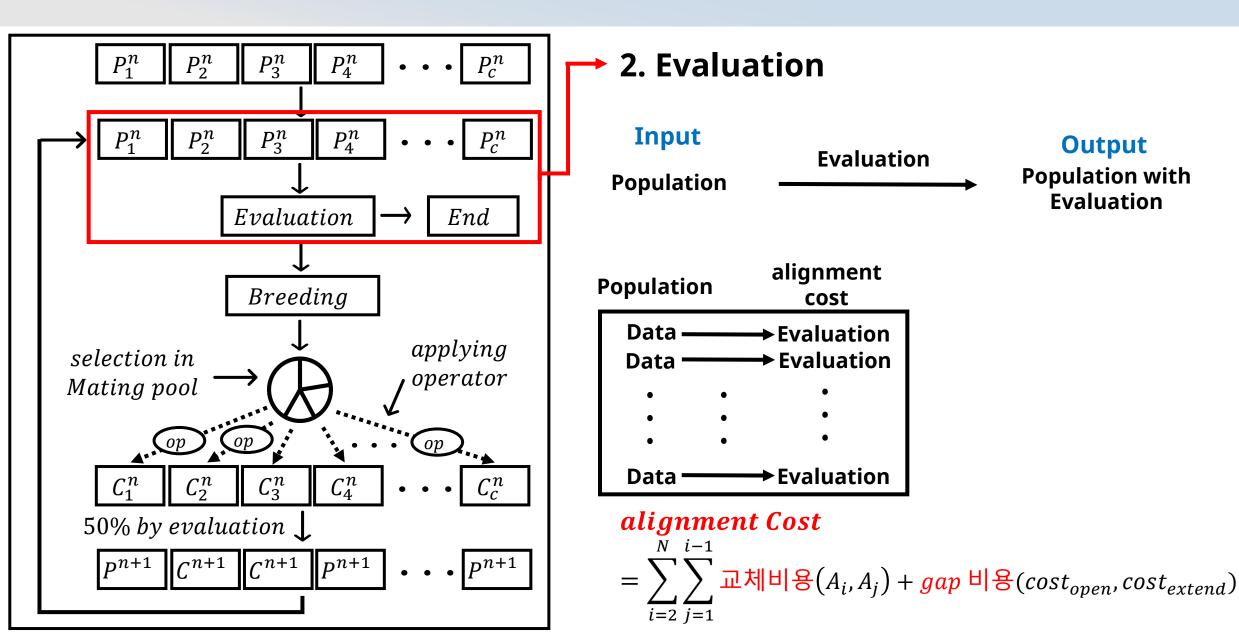
- Gap 비용은 alignment의 기호 중 '一'에 부과되는 비용이다.
- 여기서 사용된 gap 비용은 natural affine gap 비용이다.

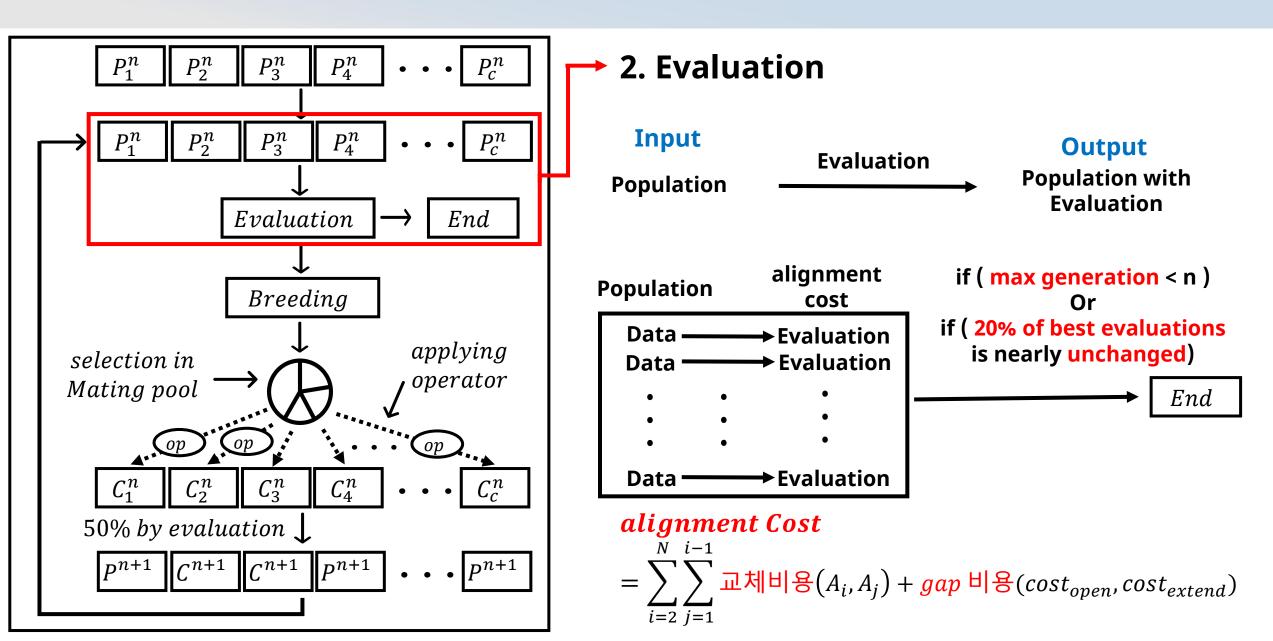
 $cost_{open}$: -1 $cost_{extend}$: -0.5라 가정

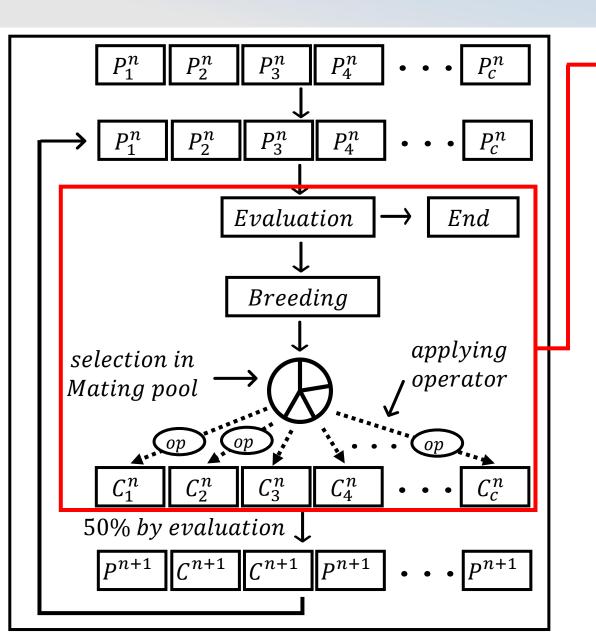
4-1.5-2+9-1.5-0.5-2+5 = 10.5

$$cost_{open} + cost_{extend} = -1 - 0.5 = -1.5$$

 $cost_{extend} = -0.5$

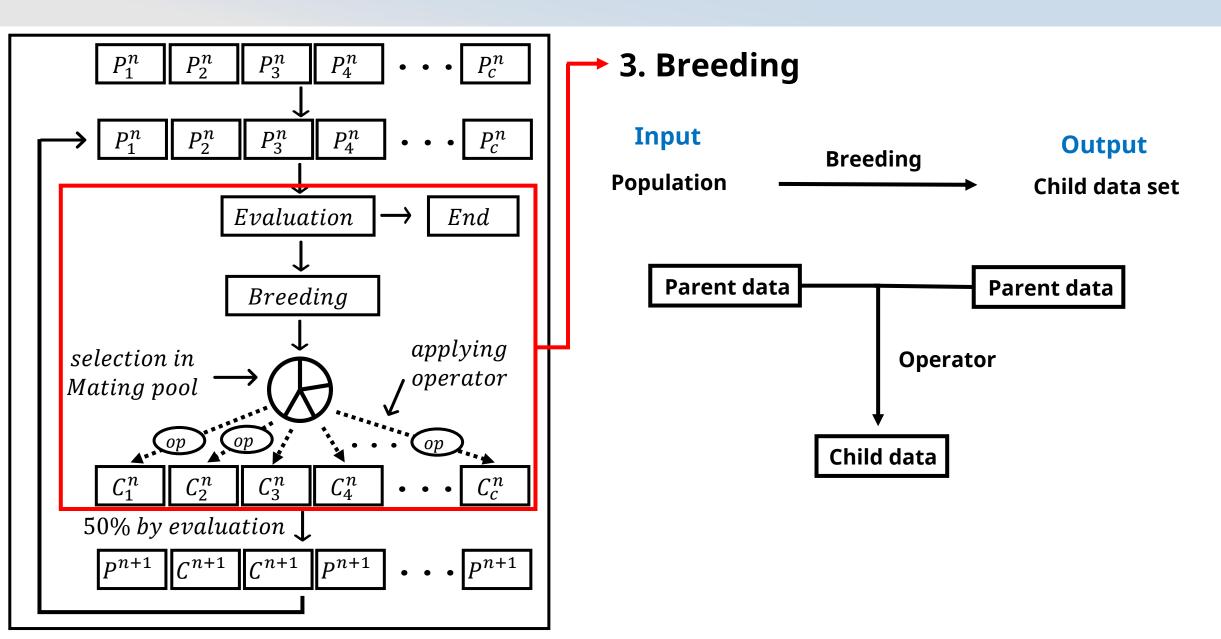


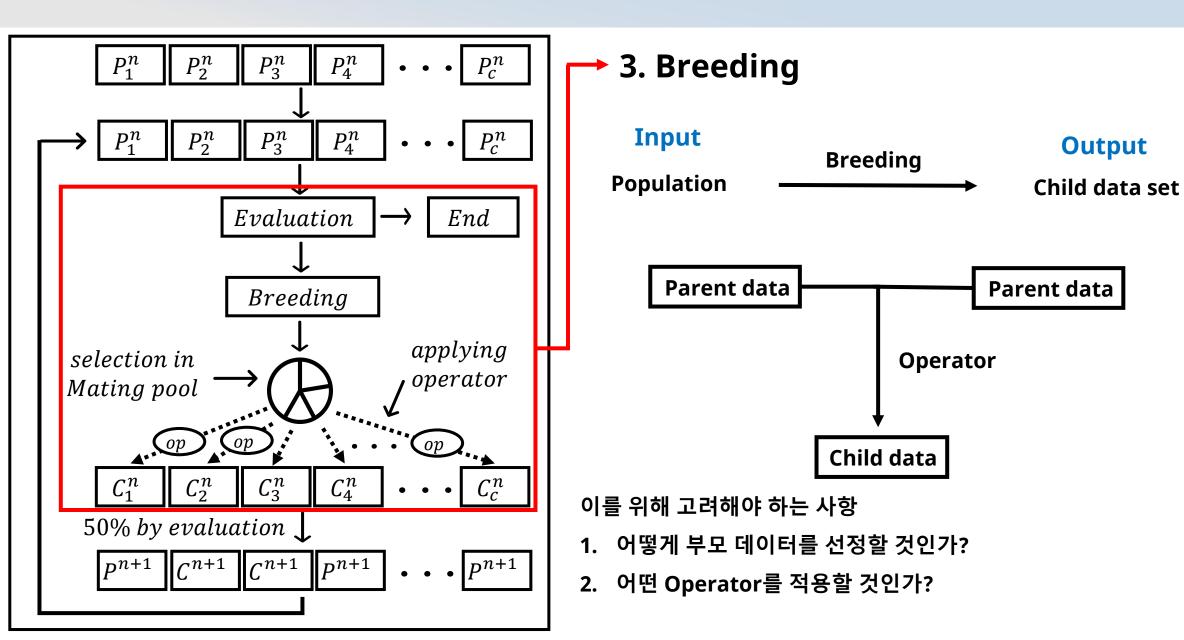


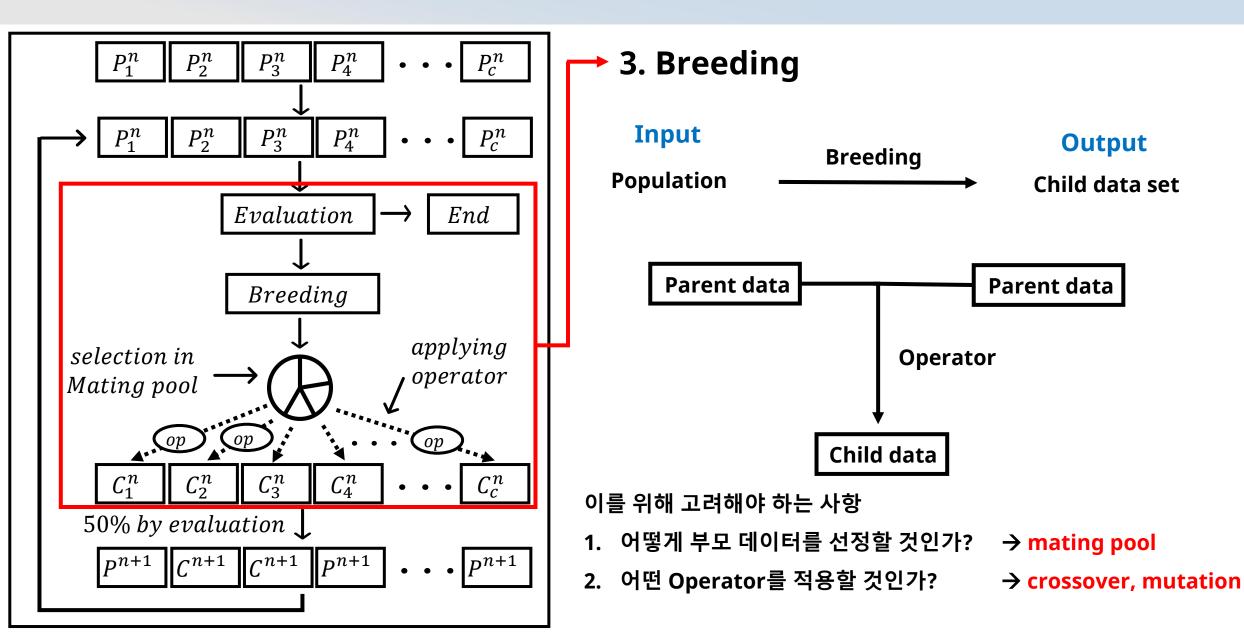


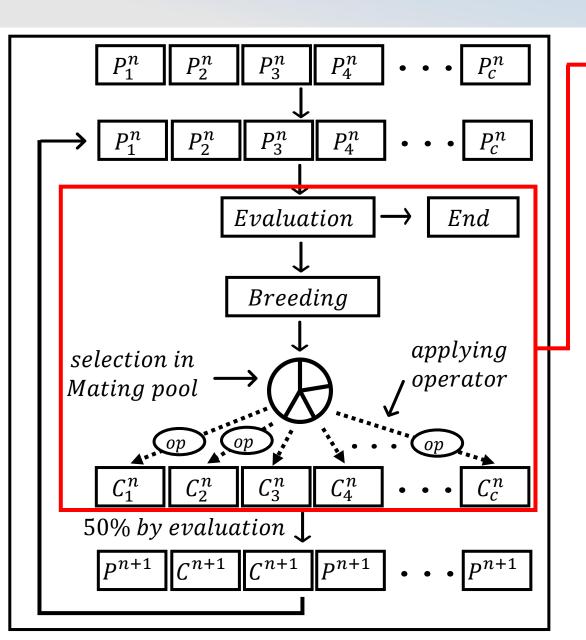
3. Breeding

Input
Breeding
Output
Child data set







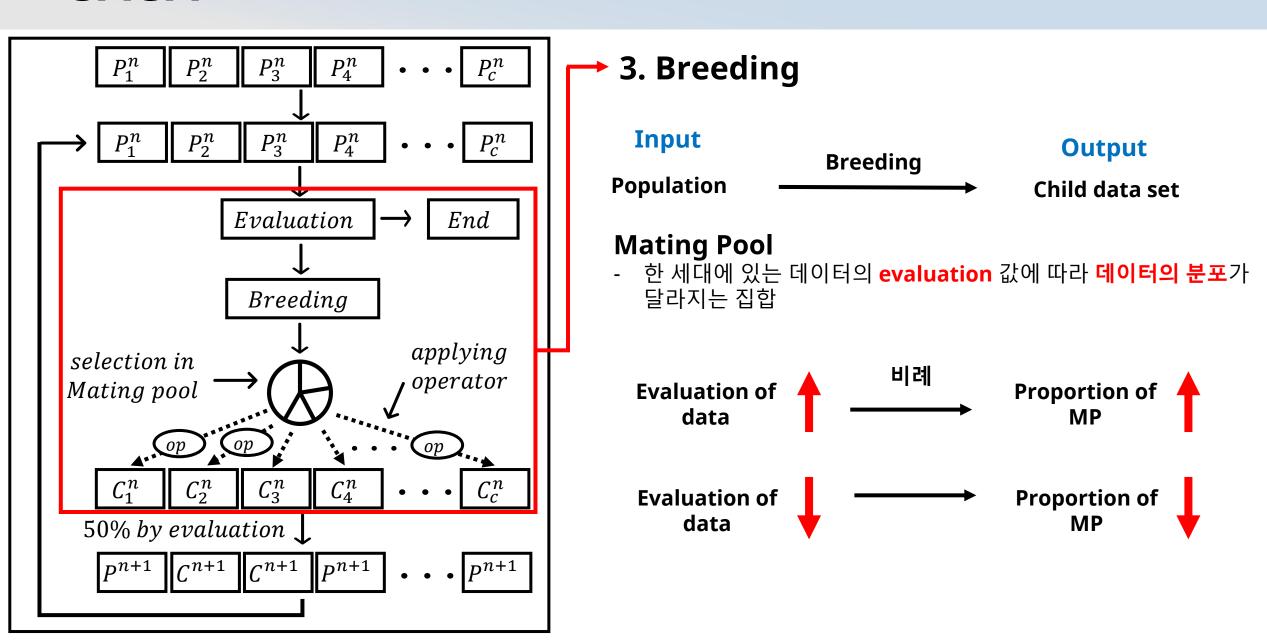


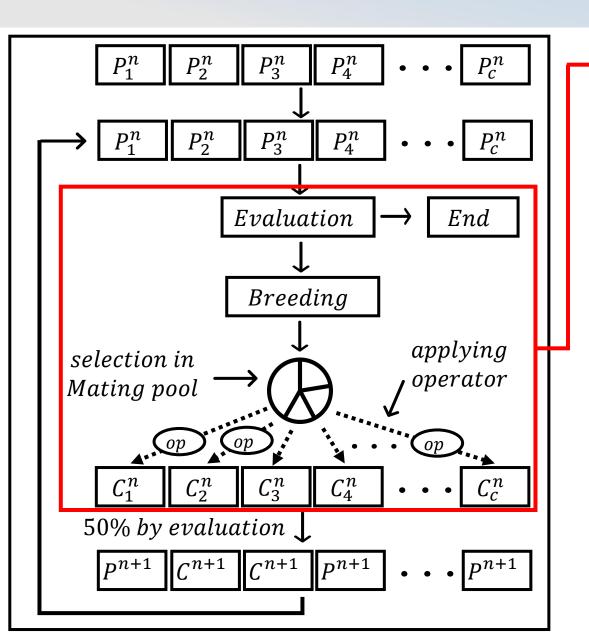
3. Breeding



Mating Pool

 한 세대에 있는 데이터의 evaluation 값에 따라 데이터의 분포가 달라지는 집합





3. Breeding

Input

Breeding

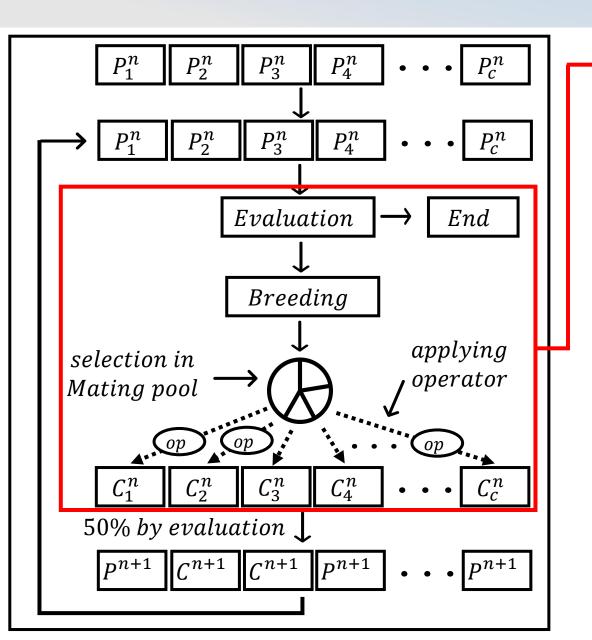
Output

Child data set

Mating Pool

- 한 세대에 있는 데이터의 evaluation 값에 따라 데이터의 분포가 달라지는 집합

MP안의 데이터 수 =
$$\left[\frac{evaluation\ of\ data}{sum\ of\ evaluations\ of\ population}\times 100\right]$$



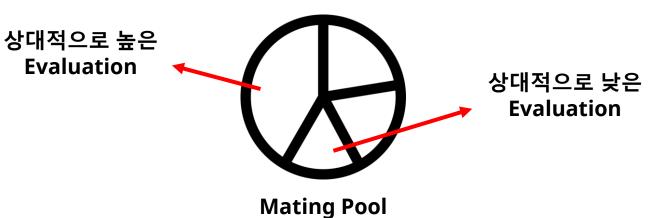
3. Breeding

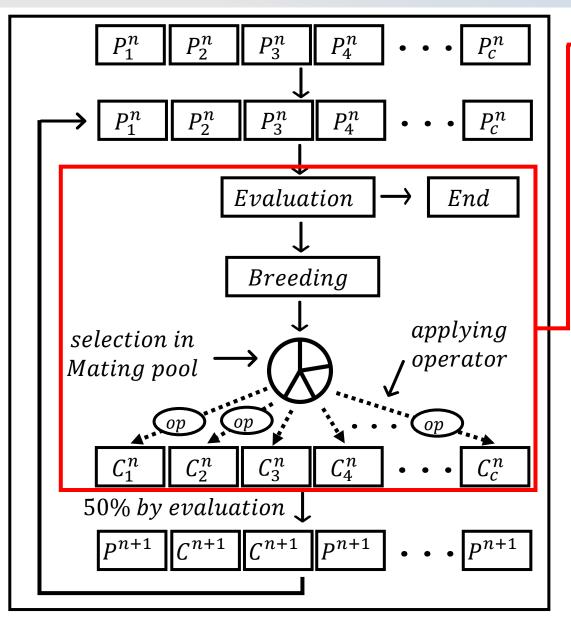


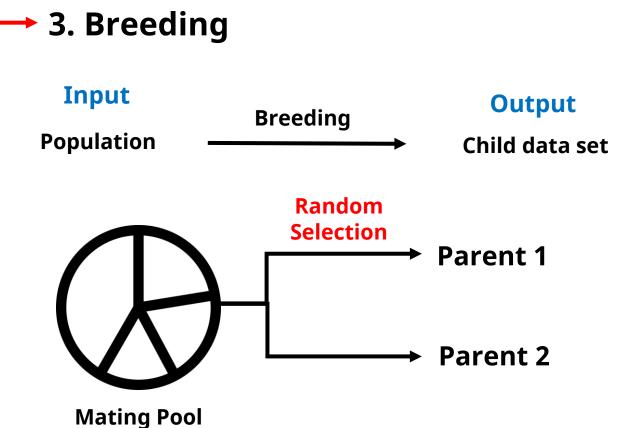
Mating Pool

- 한 세대에 있는 데이터의 evaluation 값에 따라 데이터의 분포가 달라지는 집합

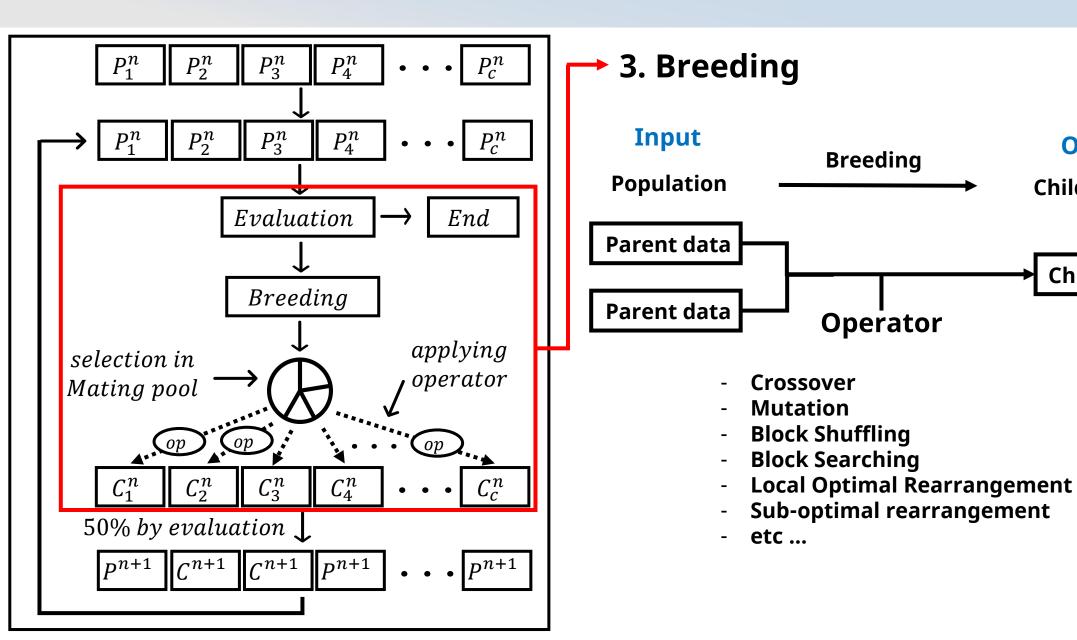
MP안의 데이터 수 =
$$\left[\frac{evaluation\ of\ data}{sum\ of\ evaluations\ of\ population} \times 100\right]$$







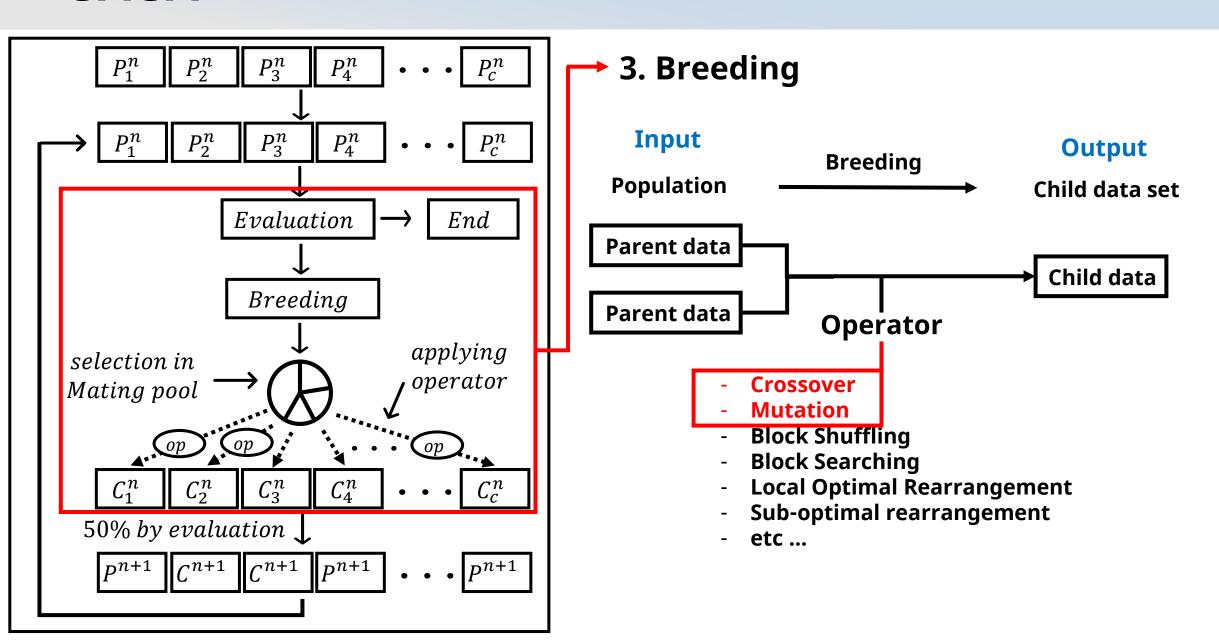
- 생성된 Mating Pool에서 랜덤으로 두 부모 데이터를 선택한다.
- 즉, Evaluation 값이 높을 수록 부모가 될 확률이 크다.
- = 좀 더 <mark>우월한 Evaluation을 가진 부모</mark>를 이용하여 <mark>새로운 세대를 구성</mark>

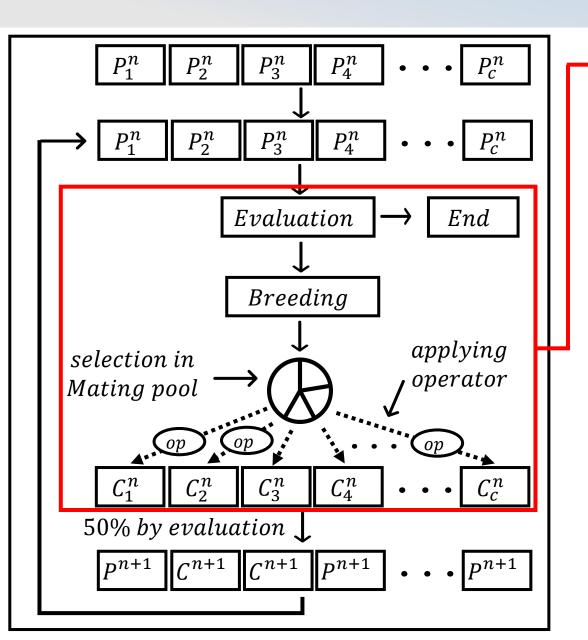


Output

Child data set

Child data

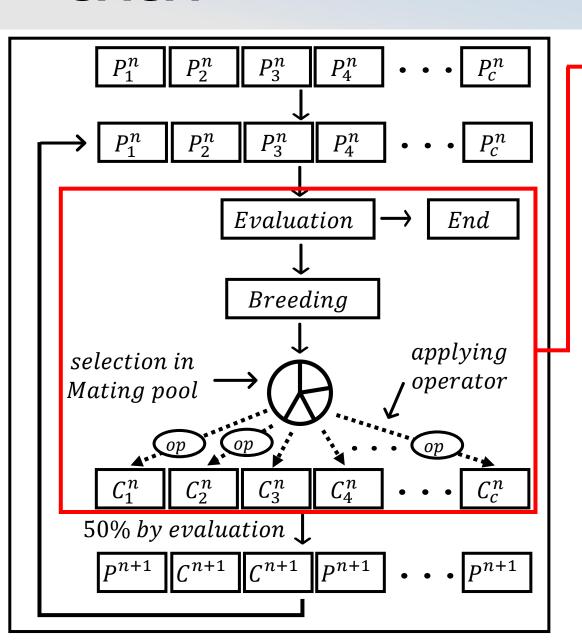




3. Breeding

Input
Breeding
Output
Child data set

- 두 부모 데이터를 이용하여 **새로운 자식 데이터를 생성**하는 연산자
- 각 데이터의 sequence들을 랜덤으로 섞는다.



3. Breeding

Input

Breeding

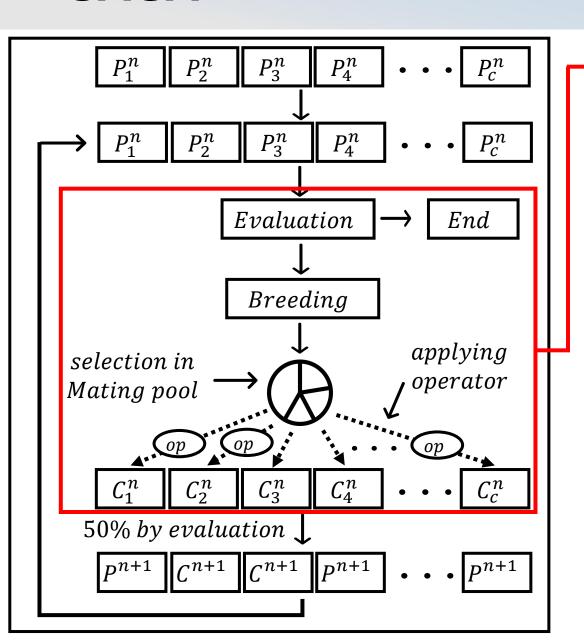
Population

Output

Child data set

- 두 부모 데이터를 이용하여 **새로운 자식 데이터를 생성**하는 연산자
- 각 데이터의 sequence들을 랜덤으로 섞는다.





3. Breeding

Input

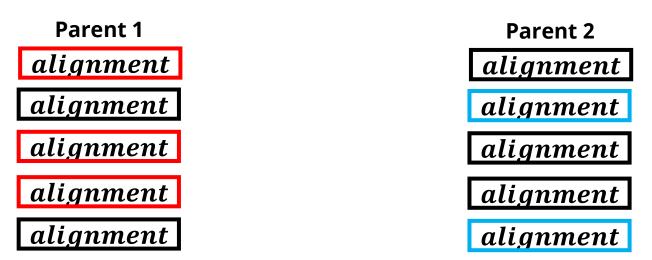
Breeding

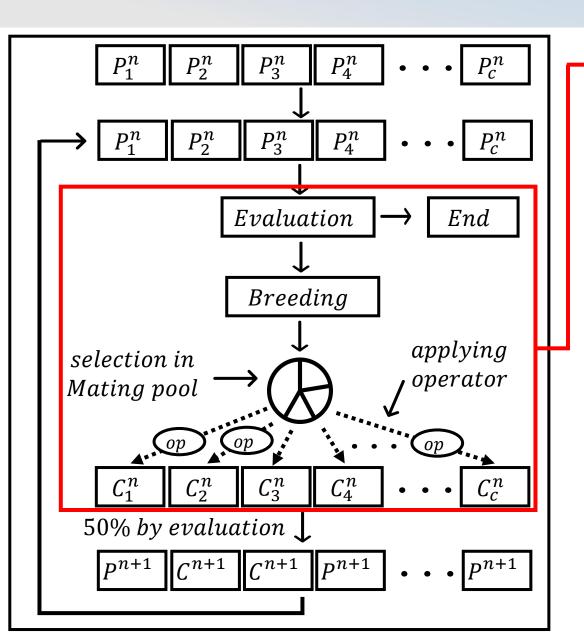
Population

Output

Child data set

- 두 부모 데이터를 이용하여 **새로운 자식 데이터를 생성**하는 연산자
- 각 데이터의 sequence들을 랜덤으로 섞는다.



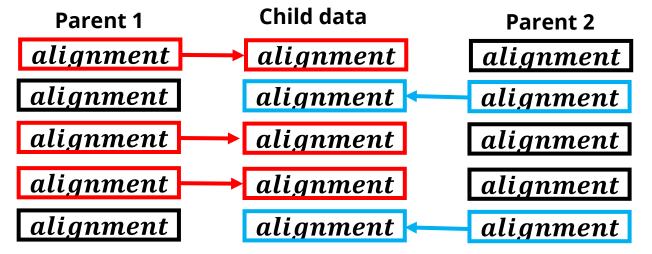


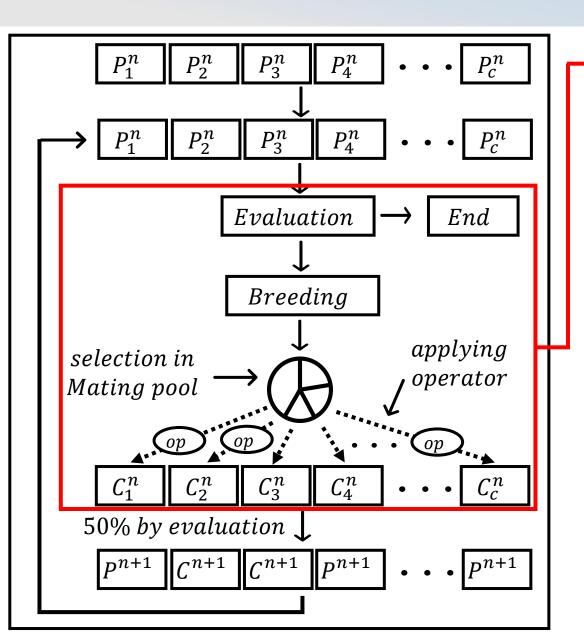
3. Breeding

Input
Breeding
Population

Child data set

- 두 부모 데이터를 이용하여 **새로운 자식 데이터를 생성**하는 연산자
- 각 데이터의 sequence들을 랜덤으로 섞는다.



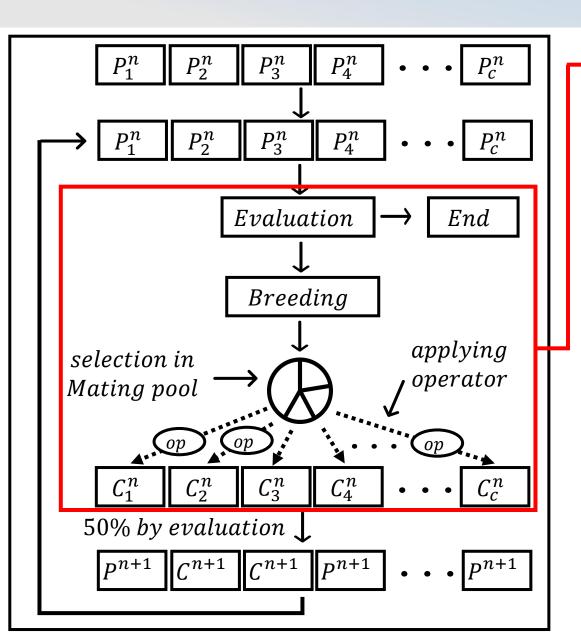


3. Breeding



Mutation

- · 하나의 데이터의 일정 부분을 변형하는 연산자
- 실제 변이가 일어날 확률은 낮기에 낮은 확률(5%)로 발생하도록 설정



3. Breeding

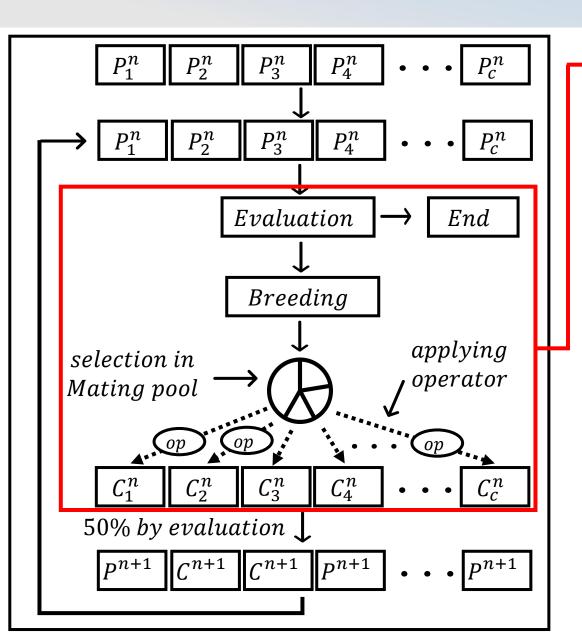


Mutation

- 하나의 데이터의 일정 부분을 변형하는 연산자
- 실제 변이가 일어날 확률은 낮기에 낮은 확률(5%)로 발생하도록 설정

1. Gaps Removal

데이터의 sequence에서 연속된 – 기호를 찾아낸 뒤 없애는 mutation



3. Breeding

Input

Breeding

Output

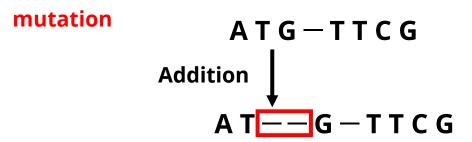
Child data set

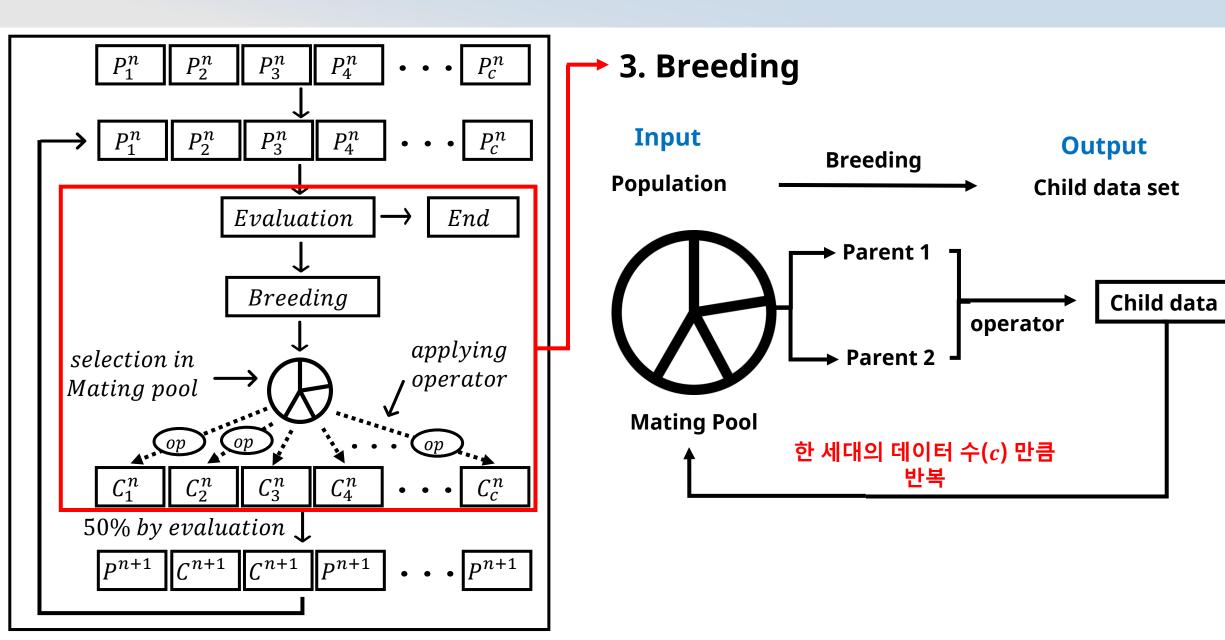
Mutation

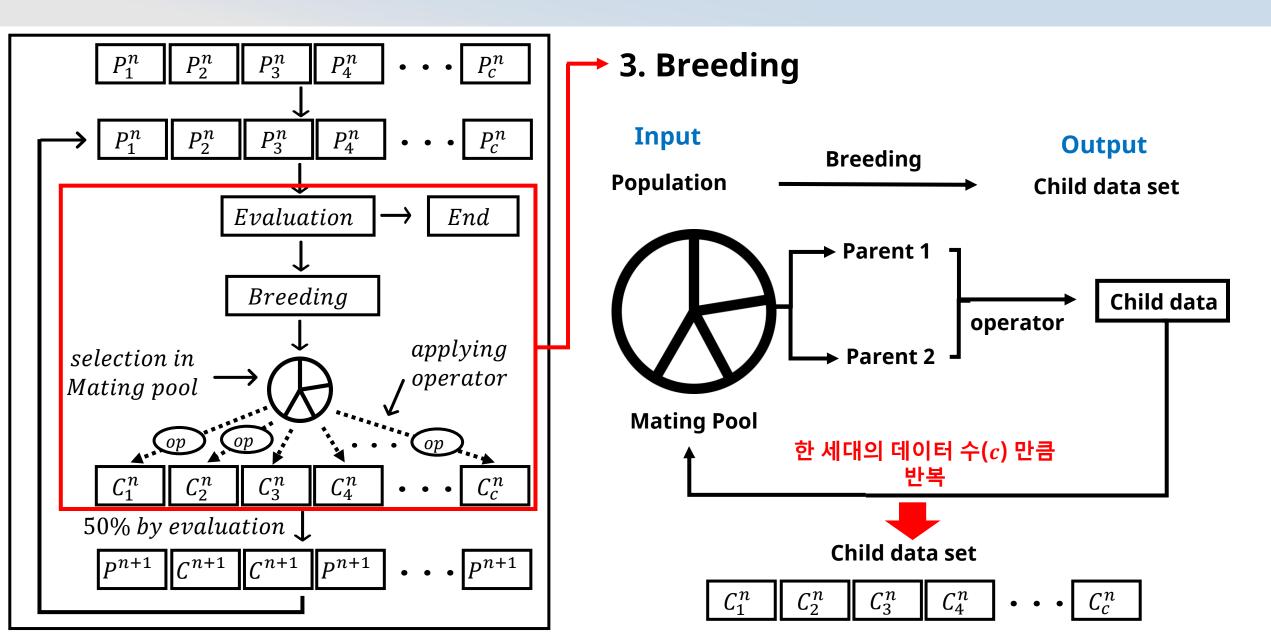
- 하나의 데이터의 일정 부분을 변형하는 연산자
- 실제 변이가 일어날 확률은 낮기에 낮은 확률(5%)로 발생하도록 설정

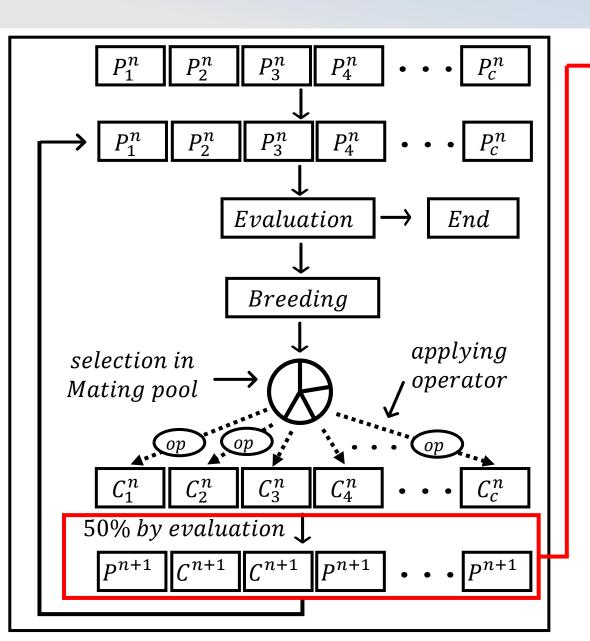
2. Gaps Addition

데이터의 sequence에 랜덤 위치에 무작위 개수의 – 기호를 추가하는



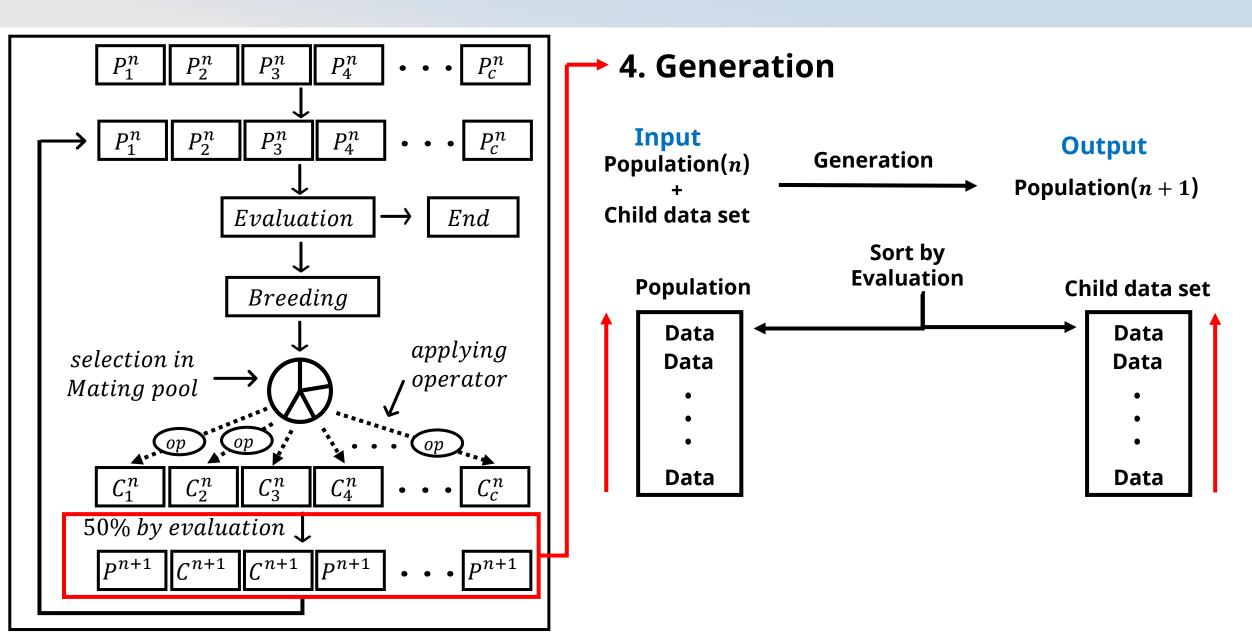


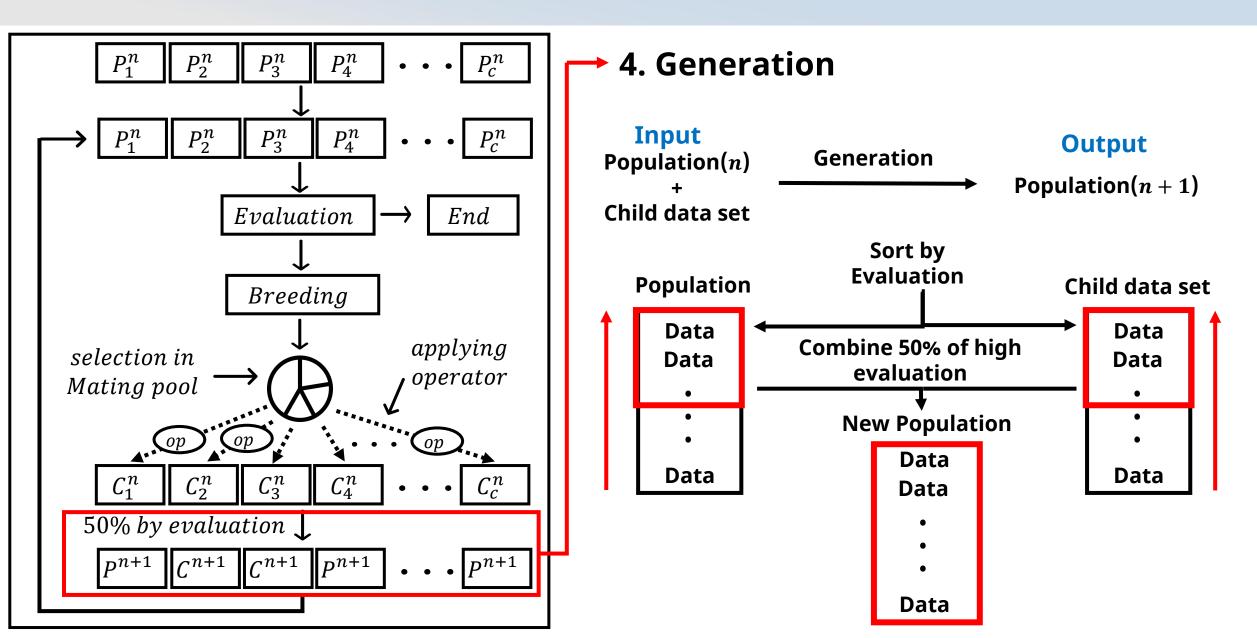


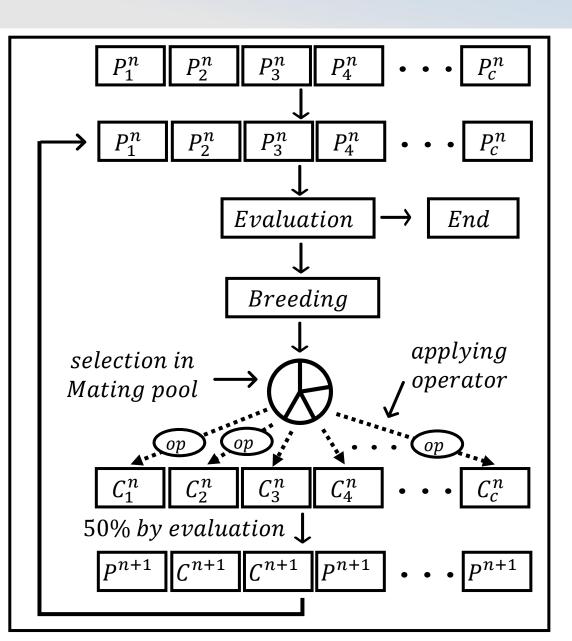


4. Generation

Input Population(n) + Generation Population(n + 1) Child data set





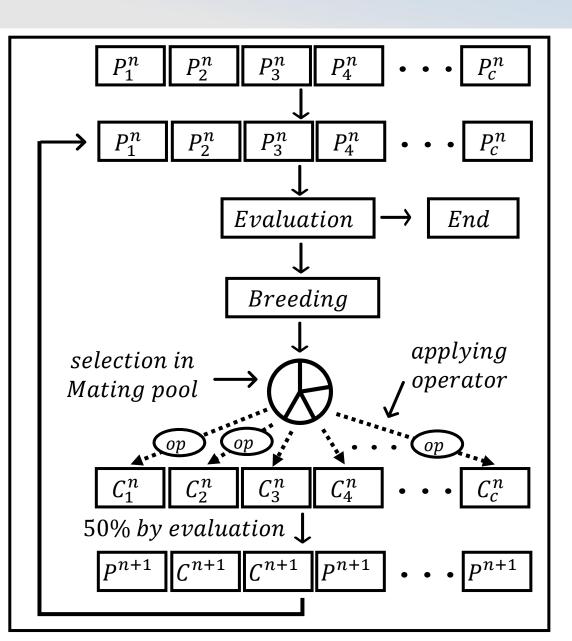


5. Iteration and End



Iteration

- 이후 2~4 과정이 반복된다.
- 반복하는 동안 각 세대 별 <mark>최대 Evaluation</mark>과 그에 대한 data 저장
- Evaluation 값이 변하지 않거나 최대 반복 수에 도달하면 종료



5. Iteration and End

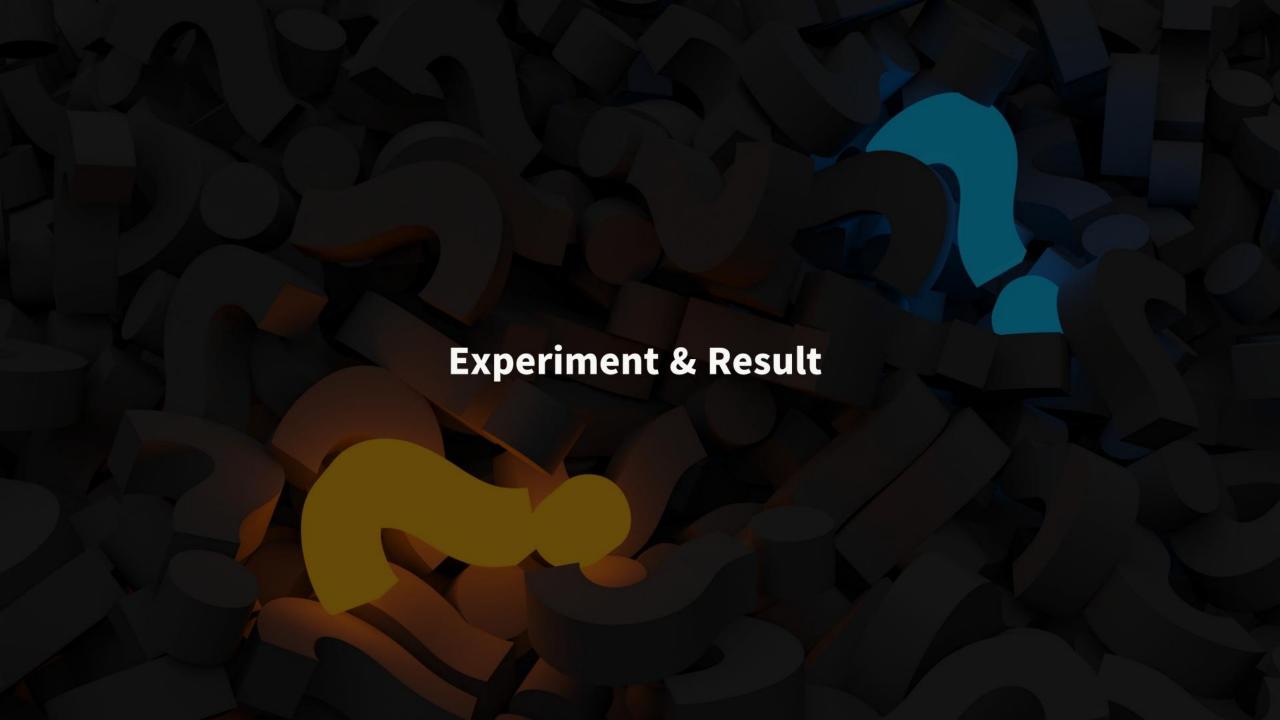


Iteration

- 이후 2~4 과정이 반복된다.
- 반복하는 동안 각 세대 별 <mark>최대 Evaluation</mark>과 그에 대한 data 저장
- Evaluation 값이 변하지 않거나 최대 반복 수에 도달하면 종료

End

- 최대 Evaluation 값과 그에 대한 **Multiple Alignment** 출력
- 동일한 Input을 다른 MSA 알고리즘인 clustal w를 통해 얻은
 Multiple Alignment 와 Evaluation 값 출력



Default Setting

- 한 번의 시행을 할 때 사용되는 변수들의 기본 값은 다음과 같다.

Setting	Explanation	Default
Max generation	최대 반복 횟수	400
Min generation	최소 반복 횟수	200
Chromosome	한 population의 총 데이터 수	100
Opening gap	gap cost에서 처음에만 주는 비용	-1
Extension gap	gap cost에서 지속적으로 주는 비용	-0.5
Mutation rate	데이터에 변이가 발생할 확률	5%

Input Format

- 입력에 사용된 예시는 다음과 같다.

Same Length

CTATCGAGTCTTCCCTCCTCCTTCTCTGCCCCCTCCGCTCCGCTGGAG
CCCTCCACCCTACAAGTGGCCTACAGGGCACAGGTGAGGCGGACTGGAC
AGCTCCTGCTTTGATCGCCGGAGATCTGCAAATTCTGCCCATGTCGGGGC
TGCAGAGCACTCCGACGTGTCCCATAGTGTTTCCAAACTTGGAAAGGGCG
GGGGAGGGCGGAGGATGCGGAGGGCGGAGGTATGCAGACAACGAGTCAG
AGTTTCCCCTTGAAAGCCTCAAAAAGTGTCCACGTCCTCAAAAAAGAATGGA
ACCAATTTAAGAAGCCAGCCCCGTGGCCACGTCCCTTCCCCCATTCGCTC
CCTCCTCTGCGCCCCCGCAGGCTCCTCCCAGCTGGCCCCC
CAGCCCCAGCCCTCCCATTGGTGGAGGCCCTTTTTGGAGGCACCCTAGGGC
CAGGGAAACTTTTGCCGTATAAATAGGGCAGATCCGGGCTTTATTATTTT

Example of Input 1

Input Format

- 입력에 사용된 예시는 다음과 같다.

Different Length

Example of Input 2

Result of Input 1

- 위의 예시 Input 1를 SAGA로 Multiple Alignment와 그에 대한 Evaluation 결과

→ 최대 Evaluation: 1212.0

SAGA Result : 1212.0 C-C-C-TC-CACCCTA-C-AAG-TGGCCT-AC-AG-G-G-CACAG-GTG----AG-G-C-G-GG-ACTGGAC--AGCT-CCT-GCT-T-TG---ATC-GC-C-GGAGA-T--CTG-CA-A--AT-T-CT-GCCCATG-TCGGGGC-T-GCA-GAGCACTCC-GACGT-G-TCCCA-T-AGTGT-T--TCCA-AA-CT-TGGA-A-A------GGGCG---G-GG-GAGG-GCG-G-G-A-G-GATGCGG-A---GG-G-CG-G-AGGTA-TGCAG-ACAACGAGTCAG-A-G-T-TT-C-CCCTT-G-AAA---GCCTCAA-AAGTGTC-CAC-GTC-CTCAA-A-A-A-A-GAA-TGGA-----A----CCA-A-TTT-AAGAAGCCA-GCCC-C--GTGG-CCACGTCC-CT-T-C-C-C-C----C-ATT-CGC-TC-C--CT----C-CTCT-G-CGC-C--CCCG-C-AG-GCTCCTCCC----AG-CTGTGGC-T-G-CCCGGGCC-CC-C----AGC-C-CC---AGC-C-CTC-CC--AT-TGGT-GGAGGCCCTT-TTGG-AGG-CA--CCC-TAGGGC CA-G-G-GAAACTT-TTGCCGT-A-T----A-A-A-T-AGGGCA-G-A-TCCG-G-GCTT-TA-TTA-T---TT

Result of Input 1

- 위의 예시를 Clustal W를 이용하여 Multiple Alignment와 그에 대한 Evaluation 결과

```
▶ 최대 Evaluation: 1856.5
clustal w Result : 1856.5
-AGTTTCCCCTTGAAAG-----AAAAG
---ACCAATTTAAGAAG------CCAGCCCCGTGGCCACGTCCCTT----CCCCC
  ----TGC--AGAGCACTCCG---ACGTGTCCCATAGTGTTTCCAAACTTGGAAAGGG
CTATCGAGTCTTCCCTCCCTCCT---TCTCTGCCCCCT-----CCG-----CTCCC
-AGCTCCTGCTTTGATCGCCGGA---GATCTGCAAATT-----CTG-----CCCAT
       -CCCTCCACCCTACAAGTGGCCTACAGGGC-----ACAGGTGAGGCGGG
  ---CCTCCTCTGCGCCCCCGCAGGCTCCTCCCAGCTGTGGCTGCCCGGG-----CCCC
        ----CAGCCCCA----GCCCTCCCATTGGTGGAGGCCCTTTTGGAGGCAC
---GGGGAGGCGGGAG-----GATGCGGAGGCGGAGGTATGCAGACAACGAG
          TTGCCG-----TATAAATAGGGCAGATCCGGGCTT--TATTAT
                      최종 Multiple Alignment
```

Result of Total Input cases

- 준비한 각 Input case에 대하여 **10번** 실행한 결과 **평균 Evaluation**과 **Clustal W의 Evaluation** 비교

Test Case	SAGA	Clustal W
1	1047.95	1856.5
2	475.7	786.5
3	602.2	883.0
4	444.3	1185.5
5	375.3	1125.5
6	695.45	1328.0
7	848.9	1179.0
8	1114.2	2229.5
9	871.1	1934.0

Result of Total Input cases

- 준비한 각 Input case에 대하여 10번 실행한 결과 <mark>평균 Evaluation</mark>과 Clustal W의 Evaluation 비교

Test Case	SAGA	Clustal W	
1	1047.95	1856.5	
2	475.7	786.5	
3	602.2	883.0	
4	444.3	1185.5	
5	375.3	1125.5	
6	695.45	1328.0	
7	848.9	1179.0	
8	1114.2	2229.5	
9	871.1	1934.0	

Result of Total Input cases (Different Gap)

- Input case 1에 대하여 5번 실행한 결과 <mark>평균 Evaluation</mark>과 Clustal W의 Evaluation 비교
- 단, Gap 비용을 다르게 설정

G	ap_{ope}	n Gap_{extend}	SAGA	Clustal W
	1	0.5	1,138.1	1856.5
	3	0.5	-332.8	1538.5
	6	0.5	-2392.9	1061.5
	1	0.05	1899.77	2342.94
	1	0.2	1448.83	2180.8
	1	0.7	633.93	1640.3

 Gap_{open} 이 커질 수록 Evaluation 또한 현저히 낮아지는 것을 확인 가능

Result of Total Input cases (Different Gap)

- Input case 1에 대하여 <mark>5번</mark> 실행한 결과 <mark>평균 Evaluation</mark>과 **Clustal W의 Evaluation** 비교
- 단, Gap 비용을 다르게 설정

Gap_{open}	Gap_{extend}	SAGA	Clustal W	
1	0.5	1,138.1	1856.5	
3	0.5	-332.8	1538.5	
6	0.5	-2392.9	1061.5	
1	0.05	1899.77	2342.94	
1	0.2	1448.83	2180.8	
1	0.7	633.93	1640.3	

이전과 같이 Clustal W를 뛰어넘지 못함

Result of Total Input cases (Different Mutation Prob)

- Input case 1에 대하여 <mark>5번</mark> 실행한 결과 <mark>평균 Evaluation</mark>과 **Clustal W의 Evaluation** 비교
- 단**, Mutation 확률**을 다르게 설정

Mutation Prob	SAGA 평균	SAGA 최대값	Clustal W
5%	1173.7	1418.0	1856.5
10%	1333.3	1491.0	1856.5
20%	1140.8	1443.0	1856.5
2.5%	942.5	1068.0	1856.5
1%	883.2	1072.5	1856.5
0.5%	677.4	812.0	1856.5

Mutation 10% 확률로 설정했을 때 상대적 좋은 성능을 보임

Result of Total Input cases (Different Mutation Prob)

- Input case 1에 대하여 <mark>5번</mark> 실행한 결과 <mark>평균 Evaluation</mark>과 **Clustal W의 Evaluation** 비교
- 단, Mutation 확률을 다르게 설정

Mutation Prob	SAGA 평균	SAGA 최대값	Clustal W
5%	1173.7	1418.0	1856.5
10%	1333.3	1491.0	1856.5
20%	1140.8	1443.0	1856.5
2.5%	942.5	1068.0	1856.5
1%	883.2	1072.5	1856.5
0.5%	677.4	812.0	1856.5

그러나 마찬가지로 Clustal W보다는 성능이 좋지 않음

Result of Total Input cases (Different Chromosome)

- Input case 1에 대하여 <mark>5번</mark> 실행한 결과 <mark>평균 Evaluation</mark>과 <mark>Clustal W의 Evaluation</mark> 비교
- 단, **Chromosome(데이터의 수)**을 다르게 설정

Chromosome	SAGA 평균	SAGA 최대값	Clustal W
100	966.5	1085.0	1856.5
200	1,151	1233.5	1856.5
300	1,250	1613.5	1856.5
400	1246.8	1457.5	1856.5

대략 300개의 데이터를 가질 경우 높은 Evaluation을 가지는 경향을 보임

Result of Total Input cases (Different Chromosome)

- Input case 1에 대하여 <mark>5번</mark> 실행한 결과 <mark>평균 Evaluation</mark>과 <mark>Clustal W의 Evaluation</mark> 비교
- 단, **Chromosome(데이터의 수)**을 다르게 설정

Chromosome	SAGA 평균	SAGA 최대값	Clustal W
100	966.5	1085.0	1856.5
200	1,151	1233.5	1856.5
300	1,250	1613.5	1856.5
400	1246.8	1457.5	1856.5

그러나 마찬가지로 Clustal W보다는 성능이 좋지 않음

Discussion about Result

- 대부분의 실행에 대하여 좋지 않은 결과가 나온 것을 확인할 수 있었다.
- 여러 방면에서 이에 대한 이유를 고찰해보았을 때 추측된 이유는 다음과 같다.

Discussion about Result

- 대부분의 실행에 대하여 좋지 않은 결과가 나온 것을 확인할 수 있었다.
- 여러 방면에서 이에 대한 이유를 고찰해보았을 때 추측된 이유는 다음과 같다.

1. 모델의 단순화

- 현재 제작된 모델은 참고한 논문에 비해 **간단화 된 모델**이다.
- 사용한 연산자의 수가 논문에 비하여 GA에서 의미적인 요소들만 사용됨

개발 : 사용한 연산자를 crossover와 mutation으로 한정

논문 : 총 **22개의 연산자**를 사용

- 모델이 반복하는 동안 연산자를 동일한 순서로 사용함

개발: crossover를 적용한 후 mutation을 확률적으로 진행함

논문: Dynamic Scheduling을 사용하여 효율적으로 연산자 사용

Discussion about Result

- 대부분의 실행에 대하여 좋지 않은 결과가 나온 것을 확인할 수 있었다.
- 여러 방면에서 이에 대한 이유를 고찰해보았을 때 추측된 이유는 다음과 같다.

2. 1990년대 논문 참조

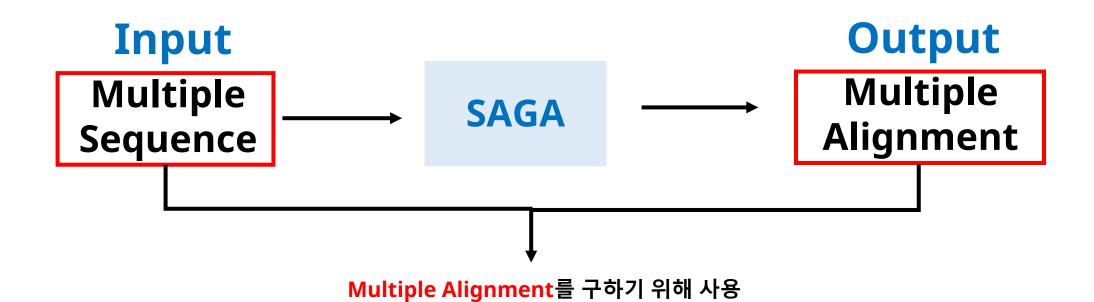
- SAGA를 만들기 위해 참조한 논문들은 각각 **1996년, 1997년대**에 발표된 논문들이다.
- 현재는 더 좋은 MSA 기법들이 개발되고 있다.
- → 따라서, Genetic Algorithm을 MSA에 사용하는 것은 적합하지 않을 수 있다.

Conclusion

- 결국 위의 실험 결과와 고찰을 통해 개발한 SAGA는 별로 효과적이지 못한 것을 알 수 있다.
- 그러나, 앞서 설명했듯 Genetic Algorithm을 MSA에 **다르게 활용할 수 있는 방법**은 존재한다.

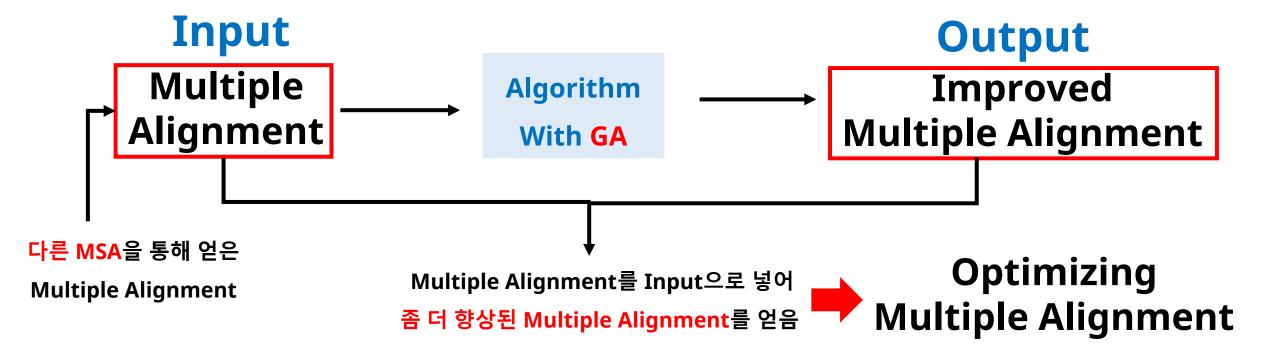
Conclusion

- 결국 위의 실험 결과와 고찰을 통해 개발한 SAGA는 별로 효과적이지 못한 것을 알 수 있다.
- 그러나, 앞서 설명했듯 Genetic Algorithm을 MSA에 **다르게 활용할 수 있는 방법**은 존재한다.



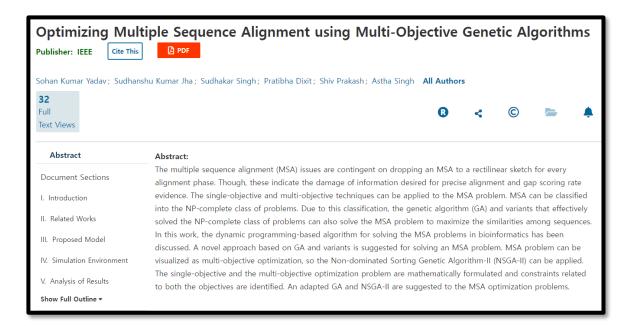
Conclusion

- 결국 위의 실험 결과와 고찰을 통해 개발한 SAGA는 별로 효과적이지 못한 것을 알 수 있다.
- 그러나, 앞서 설명했듯 Genetic Algorithm을 MSA에 **다르게 활용할 수 있는 방법**은 존재한다.

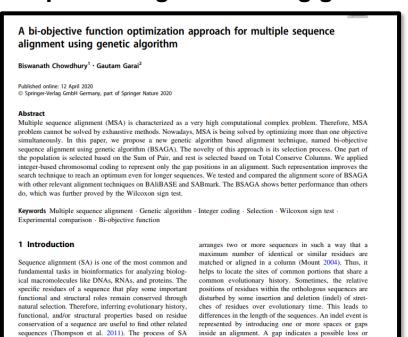


Conclusion

- 결국 위의 실험 결과와 고찰을 통해 개발한 SAGA는 별로 효과적이지 못한 것을 알 수 있다.
- 그러나, 앞서 설명했듯 Genetic Algorithm을 MSA에 **다르게 활용할 수 있는 방법**은 존재한다.
- 1. Optimizing Multiple Sequence Alignment using Multi-Objective Genetic Algorithms



2. A bi-objective function optimization approach for multiple sequence alignment using genetic algorithm



absence of a residue in a sequence with respect to other set



Thank you.



Reference

- 1. 웹으로 multi sequence 할 수 있는 곳: https://www.ebi.ac.uk/Tools/msa/clustalo/
- 2. 진화 알고리즘을 사용한 복수 염기서열 정렬

The Korean Journal of Microbiology, Vol. 35, No. 2. June 1999, p. 115-120 Copyright@1997, The Microbiological Society of Korea

진화 알고리즘을 사용한 복수 염기서열 정렬

김 진¹ · 송민동² · 최홍식³ · 장연아³ 건국대학교 자연과학대학 전산과학과1*, 분자생물학과2, 한림대학교 컴퓨터공학부3

3개 이상의 DNA 혹은 단백질의 염기서열을 정렬하는 복수 염기서열 정렬(multiple sequence alignment)은 염기서열들 사 이의 진화관계, gene regulation, 단백질의 구조와 기능에 관한 연구에 필수적인 도구이다. 복수 염기서열 정렬을 얻기 위한 기 존의 방법은 progressive pairwise alignment와 같이 빠른 실행시간 내에 만족할 만한 복수 염기서열 정렬을 제공하는 방법과, 최저의 복수 위기서열 정렬을 제공하나 실행시간이 상대적으로 긴 dynamic programming과 같은 방법 등이 있다. 본 논문에 서는 진화 알고리즘을 사용하여 기존의 방법에서 제공하는 복수 염기서열 점렬을 짧은 시간 내에 보다 개선된 복수 염기서열 점렬을 획득하게 하는 방법을 제시하였으며, 진화 알고리즘의 구성내용을 설명하였으며, 실제의 염기서염을 사용하여 이 방법 의 장점을 보였다.

KEY WORDS

multiple sequence alignment, genetic algorithm, dynamic programming, sequence comparison

생물학 역사상 가장 중요한 프로젝트의 하나인 Human Genome Project의 기본적인 목표는 인체의 게놈과 생명체의 유전자 의 염기서열의 인식을 목표로 하고 있다. 이 프로젝트에 의해 발생되는 엄청난 양의 염기서열 관련 데이터는 의약과 생물학 분야에 절대적인 영향력을 미치고 있으며 이러한 추세는 더욱 심화될 것이라 예상된다. 이러한 염기서열 관련 데이터를 처리 하여 중요한 생물학적 정보를 얻기 위해서는 전산학의 도움이 필수적이다. 전산학에서 염기서열은 스트링으로 간주된다. 본 논문에서는 게놈 프로젝트에서 파생된 가장 중요한 문제 중에 하나인 복수 염기서열 정렬(multiple sequence alignment) 문제 에 대하여 논한다(3-5, 7, 20).

이 있다. 복수 염기서열 정렬에 사용되는 또 다른 방법은 simulated annealing(11-13), 진화 알고리즘(genetic algorithm) 등

이 있다(13). Progressive pairwise alignment에 의한 방법은 짧 은 실행시간을 사용하여 만족할 만한 복수 염기서열 정렬을 생 산한다. 그러나 이 방법은 획득되는 복수 염기서열 정렬이 최적 비용을 가짐을 보장해주지 못한다 반면에 dynamic programming을 사용하여 복수연기서열 정렬을 얻는 방법은 최적, 혹은 최 적에 가까운 복수 염기서열 정렬을 제공하나, 정렬하려 하는 염 기서열의 개수가 증가함에 따라 필요한 실행시가도 지수 함수적 으로 증가하기 때문에 4~7개 이상의 염기 서열정렬에는 효율적 이지 못하며, 특정한 측정 단위에는 사용될 수 없는 등의 단점

3. SAGA: Sequence Alignment by Genetic Algorithm

© 1996 Oxford University Press

Nucleic Acids Research, 1996, Vol. 24, No. 8 1515-1524

SAGA: sequence alignment by genetic algorithm

Cédric Notredame* and Desmond G. Higgins

EMBL outstation, The European Bioinformatics Institute, Hinxton Hall, Hinxton, Cambridge CB10 1RQ, UK

Received December 5, 1995: Revised and Accepted March 4, 1996

ABSTRACT

We describe a new approach to multiple sequence alignment using genetic algorithms and an associated software package called SAGA. The method involves evolving a population of alignments in a quasi evolutionary manner and gradually improving the fitness of the population as measured by an objective function which measures multiple alignment quality. SAGA uses an automatic scheduling scheme to control the usage of 22 different operators for combining alignments or mutating them between generations. When used to optimise the well known sums of pairs objective function, SAGA performs better than some of the widely used alternative packages. This is seen with respect to the ability to achieve an optimal solution and with regard to the accuracy of alignment by comparison with reference alignments based on sequences of known tertiary structure. The general attraction of the approach is the ability to optimise any objective function that one can invent.

INTRODUCTION

The simultaneous alignment of many nucleic acid or amino acid sequences is one of the most commonly used techniques in sequence analysis. Multiple alignments are used to help predict the secondary or tertiary structure of new sequences; to help demonstrate homology between new sequences and existing families; to help find diagnostic patterns for families; to suggest

There are two main alternatives to progressive alignment. One is to use hidden Markov models (HMMs; 5) which attempt to simultaneously find an alignment and a probability model of substitutions, insertions and deletions which is most self consistent. Currently, this approach is limited, in practice, to cases with very many sequences (e.g. 100 or more) but does have the great advantage of a sound link with probability analysis. A second approach is to use objective functions (OFs) which measure multiple alignment quality and to find the best scoring alignment. If the OF is well chosen or is an accurate measure of quality, then this approach has the advantage that one can be confident that the resulting alignment really is the best by some criterion. Unfortunately, the number of possible alignments which must be scored in order to choose the best one becomes astronomical for more than four or five sequences of reasonable length.

Two solutions to this problem exist. The MSA program (6,7) attempts to narrow down the solution space to a relatively small area where the best alignment is likely to be. It then guarantees finding the best alignment in this reduced space. Even with this reduction, it is limited to small examples of around seven or eight sequences at most. Nonetheless, it is the only method we know of that seems capable of finding the globally optimal alignment or close to it, starting with completely unaligned sequences. A second approach is to use stochastic optimisation methods such as simulated annealing (8), Gibbs sampling (9) or genetic algorithms (GAs; 10). Simulated annealing has been used on numerous occasions for multiple alignment (e.g. 11-13) but can be very slow and usually only works well as an alignment improver i.e. when the method is given an alignment that is already close to optimal and is not trapped in a local minimum. Gibbs sampling has been very successfully applied to the problem

CTATCGAGTCTTCCCTCCCTCCTTCTCTGCCCCCCTCCGCTCCCGCTGGAG CCCTCCACCCTACAAGTGGCCTACAGGGCACAGGTGAGGCGGGACTGGAC AGCTCCTGCTTTGATCGCCGGAGATCTGCAAATTCTGCCCATGTCGGGGC TGCAGAGCACTCCGACGTGTCCCATAGTGTTTCCAAACTTGGAAAGGGCG GGGGAGGCGGAGGATGCGGAGGCGGAGGTATGCAGACAACGAGTCAG AGTTTCCCCCTTGAAAGCCTCAAAAGTGTCCACGTCCTCAAAAAGAATGGA ACCAATTTAAGAAGCCAGCCCCGTGGCCACGTCCCTTCCCCCATTCGCTC CCTCCTCTGCGCCCCGCAGGCTCCTCCCAGCTGTGGCTGCCCGGGCCCC CAGCCCCAGCCCTCCCATTGGTGGAGGCCCTTTTTGGAGGCACCCTAGGGC CAGGGAAACTTTTGCCGTATAAATAGGGCAGATCCGGGCTTTATTAT

SAGA Result: 1578.5

```
CTAT-CGA-GTCT-T-CCCTCCC-T-CC-T-TC-T-C-T-GCCCC-C-TC-C-GC-TC-CCGCT-GGAG--
----AG--CTCCTG-CT-TTGATCGCCG-G-AGA--TC-TG-CAAATTCTGC-CC-ATGTCGGG-GC----
T-GCA-GAGCACTCC-GACGT-G-TCCCA-T-AGTGT-T--TCCA-AA-CT-TGGA-A-A----GGGCG---
G-G-G-G-A-G-G-G-CG-G-GAG-GAT-GCG-GAGGGCG-G-AGGT-AT-GCAGACA-A-CGAGTCAG---
-AGT-T-T-C-C-CC-TTGA-A-AGC-C-TCAAA-AGT-GTCCACG-TCCT---CA-AA-AA-GAA-TGGA
--AC-CA-A-T-TTA-A-G-A-A-GCCA-GCC--C-CGTGG-CCA-CGTCC-CTTCCCC-CATT-CGC-TC-
C--CT----C-CTCT-G-CGC-C--CCCG-C-AG-GC-TCCTCCC-AG-CTGTGGC-T-G-CCCGGGCCCC-
C----AGC-C-CC---AGC-C-CTCCC--AT-TGGT-GGAGGCCCTT-TTGGAGG-CAC-CCTAGGGC----
CAG-G-GAAACTTTT--GC-C-GT--ATA-A-A-T-AG-GG-CA-GATCCGG-G-CT-T-TATTA-T-T-TT
clustal w Result : 1856.5
-AGTTTCCCCTTGAAAG------CCTCAAAAGTGTCCACGTCCTCA----AAAAG
---ACCAATTTAAGAAG-----CCAGCCCCGTGGCCACGTCCCTT----CCCCC
-----TGC--AGAGCACTCCG---ACGTGTCCCATAGTGTTTCCAAACTTGGAAAGGG
CTATCGAGTCTTCCCTCCCT---TCTCTGCCCCCCT-----CCG-----CTCCC
-AGCTCCTGCTTTGATCGCCGGA---GATCTGCAAATT-----CTG-----CCCAT
 -----CCCTCCACCCTACAAGTGGCCTACAGGGC-----ACAGGTGAGGCGGG
----CCTCCTCTGCGCCCCCGCAGGCTCCTCCCAGCTGTGGCTGCCCGGG-----CCCC
       -----CAGCCCCA----GCCCTCCCATTGGTGGAGGCCCTTTTGGAGGCAC
---GGGGAGGCGGGAG-----GATGCGGAGGCGGAGGTATGCAGACAACGAG
CAGGGAAACTTTTGCCG-----TATAAATAGGGCAGATCCGGGCTT--TATTAT
```

```
SAGA Result : 772.5
TT-T-AGCA-AT-CT--TC-T-TGC-TTCAC-T--CTA-TCA--GA-GGT-A-T-GGA--CCATC-A-CT-CT-GGT-CAA-C
T-GGTTT-TTTTTTT
T-T--TCT--CC-C-CAGA--CTGGAGCAATC-T-CT-T--ATAGTG--CAGGT-T-G-GT-TTTGAG-CTC-GA-GGCAATA-CTC-CC-GTCCT---
A-C-CTCA-GCCTCTCA-ATGCTGG-G--ATG-A-CA-A-GATA-TC-CCAGGC-A-A-GC-TTTGAA-CTT-GC-GGCAATT-CTG-CT----TT--
AA--CCT-C-CT-T--AG-TGC-TC--TCTACC-ATGAA-TCT-ATGG-GA-AGA-A-GA-AA-T-A-A-T-G-GGGGCGGG-G-G-GGGA-AACAACC
T--GCA-GCGA-G-CG-AT-GATG-A-TCACGTGACT-AGTC-C-TGCGG-G-GCGGAGGCCAT-----GT-T-GCGGG-GCA-CCCA--CGT-GA
G-GGCC-GC-ACGTCC-AC-GATC-AGTCACGTGACC--GTG-G-TGCGC-C-GC--A-GCC-----GC-C-G-GGGCGCACCCG-GCGA-GAGG---
C-A-G-C-G------GCAG-TG------
clustal w Result : 786.5
AACCTCCTTAGTGCTCTCTACCATGAATCTATGGGAAGAAGAAATAATGGGGGCG-----
-----TTTCTCCCCAGACTGGAGCAATCTCTTATAGT-GC----AGGTTGGTTTT
----TTT------AGCAATCTTCTTGCTTCACTCTATCAGAGGTATG
--GGGCCGCACGTCCACGATCAGTCACGTGACCGTGGTGCGCCGCA--GCCGCCGGGGGCG
----TGCAGCGAGCGATGATGATCACGTGACTAGTCCTGCGGGGC---GGAGGCCATGT-
-AACTCTGGGCATCAGTTCGG-----ATTAAGGTCGA-TCCGCGCA--TGCGTTCATTTA
-ACCTCAGCCTCTCAATGCTG---GGATGACAAGATATCCCAGGCA--AGCTTT-----
```

```
SAGA Result : 598.0
--G--CA---GGCTGC-CT-T-TG-GT-G-ACTCACC-GGG-TGA-ACGG-GG-GCATTGCG-AGGCAT-CCCC-T-CC-CTGG-GTTTGG---CTC--
CTGCCCACGGGGCTGA-CA-G-TA-GA-A-A-TCACA-GGC-TGT-GA-GA----CA--GCT--GG-AG-CCCA-G-CT-CTGC---TTGA-ACCTA-T
T-T-T-AGG--TCT---CT-GAT--C-CCCGCT-TCC-TC-T-TT-AGA-C-T-CCCCTA-GAGCTCAG-C-CAGTGCT-CAA-CCT-GAGG-CTGG--
G-G-GTCTCT-G-AG-G-A-AG-AG-TGAG----TTGGA-GC-TG-AGGG-G-T-CTGGG-GCTGT-CC-CCTGAG-AGAGGG-GCCAG-AGGCA----
G-T-G-TCA--AGAGCCG-GGCAG-TC-TG-A--TTGT-GGC-TCA-C-CCTC-CA-T-CA-CTC-CCAGGGCC-C-C-TGG---CCCAG-C--AGCCG
C-A-GCTCC-CA-A-CCA---CAA-TA-TC-C-TTTGG-GGT-TTG-G-CCTA-CG-G--AGCTG----GGGCG-G-A-TGA-CCCCCAAA-TAGCCCT
G-G-C-AGA-TTCC-CCCTAGAC-CCGCCCGCA--CCATG-G-TC-AGG-CATGCCCCTC----CTCAT-CGCT-GGCA-CAGCCC--AGA-----
clustal w Result : 883.0
GGCAGATTCC-----CCCTAGACCCGCCCGCACCATGGTCAGGCATGCCCCTCCT
---CAG-----CT--CCCAACCA----CAATATCCTTTGGGGTTTGGCCTACGGAG
-----GT--GTCAAGAGCCGGGCAGTCTGATTGTGGCTCACCCTCCATCA
------GGGT--CTCTGAGGAAGAGTGAGTTGGAGCTGAGGGGTCT-GGGGCT
-----TTTAGGTCTCTGATCCC----CGCTTCCTCTTTAGACTCCCCT-AGAGCT
-GCAGGCTGCCTTTGGTGACTCACCGGGTGAACGGGGGCATTGCGAGGCATCC-----
     ·CTGCCCACGGGG-CTGACAGTAGAAAT-CACAGGCTGTGAGACAGCT-GGAGCC
  ------
```

```
SAGA Result : 491.5
--AC-TT-CT-T---TTCTG-G-TGACAG-GAAA-AGGACTGGGAAGAGGGT-T-T-TCT-TC-GGGGTGC-T-TC-A-GGA-A-G-GC-AG-AAT-
ACTGTGCCTG-AGCC-AGCTGAGC-C-AGCTGA-GCCAGTT-AGC-AAG-T-G-CAT-G-GAC-T-GC-C-G-GAT-TC-GAT-T-----GGATT
G-GAT-TGCG-G--ATTGCC-AA-A--ACAAG-GAC-AGCC--ACTGGGACAGGGA--CAGC---GGGAC-AGGGA-CA-GT-GAC-G-G-CGGA--
--G-GGCAGC-AG-G--GTGGG-GGT-GGG-GT-GTCT-TTC-T-CG--ACT-GGTGAG-GAAA-AGGGA-GGGGA-GAAGA-----CTTTTTAAA
TGAG-TG-CA-A-AATTCTA-C--G-CAC--AAAGA-GAG-GGGAG-GGAGGGA-C-C--CA-GA-GGGG-GA-AAGGG-G-TGG-AGAAA-AGA--
AAGAGG-AGG-AA-GAAGTGGGAGGA-GGGAGG-G-CG-ATGG--GATACC-GGTT-AA--AAAGAGGGT-GGGCGGG-AGC-------AAC
TAACAGGCTCA-----
clustal w Result : 1185.5
AAGAGGAGGAAGAAGTGGGAGGAGGGAGGGCGATGGGATACCGGTTAAAAAG------
------GGGCAGCAGGGTGGGGTGTGTCTTTCTCGACTGGTGAGGAA
ACTTCTTTTCTGGTGACAGGAAAAGGACTGGGAAGAGGGT-TTTCTTCGGGGTGC---TT
-----ACTGTGCCTGAGCCAGCTGAGCCAGCTGAGC-CAGTTAGCAAGTGCATGGA
--GGATTGCGGATTGCCAAAACAAGGACAGCCACTGGGAC-AGGGACAGCGGGACAGGGA
-----AGGCTCA------
```

```
SAGA Result : 330.5
CTT-T-T--CTGGTGA-CA--GAA-A-AAG-G-A-CTG-G-A-A-A-A-A-A-GAGATTT-TCCTC--AGGG--TGC-T-TC-A-GGAA-AGCAGAAT
G-T-A-GC-CTC-A--GCA--GC-AGAGAGATG-G-G-CA-G-A-A-A-A-A-A-GA----CCCA-GGAA-GGC-TGTGC-CTGA-GTC-AGC-TGAGACAGC-
-TA-G--CA-A-GTG-CAT-----GGACTGC-C-AGAT-TGCCGATTG-TGG-ATC-GC-C-A--CG-A--CAA-GAACAGCCACTTGGACAG-GG-A---
C-AGC-G-G-AG-GGCAGC-AG-G-GTGGGG-GGT-GGG-GTGTCT-T-TC-T-CG-ACT-GGTG-AGAAGAAGGGAGGAGA-GAGGA-----C---
clustal w Result : 1112.5
CTTTTCTGGTGACAGAAAAGGACTGGAAAAAGAGATTTTCCTCAGGGTGCTTCAGGA--
CTAAT-----GGCTGTGGCTTGGGAGAATTACCTTTTCCTGGGTTCCCTGGAGGGTT
----CAGCGG--AGGGCAGCAGGGTGGGGGGGTGTCTTTCTCGACTGGTGAGAAG-
---GTAGCCTCAGCAGCAGAGAGATGGGCAGAAAAAGACCCAGGAAGGCTGTGCCTGAG-
  ----TAGCAA---GTGCATGGACTGCCAGATTGCCGATTGTGGATCGCCACGACAAGAAC
GAGAGAGAAAGAGGAAGAAGTGGGAGGAGGAGGGCGATGGGATACCGGTTAAAAA-
   -----GGAGCAACCAG------
```

```
SAGA Result : 559.0
TGG-G-CG-CG-GGT-CG-GCGGCC-GCGACCCG-GGAGCGGG-T-TTGCTCAGGA-A-A-AG-GCC-C--GT-CGC-CC-C-C-CAAA-C-C-C-------
C-TC-CC-TG-G----C-CG-G-C---TC-C-CT-G-CCG-GG---CCTCC-CG-G-GCCTG-G-TG-CTA-GGGCA-CCG-CG-GG-GAG-CG-CCGAATGGGA
-CGCACTG-CA-G-GGG-C-GCCA--G-AT-T-T-GGCGGGA-GGG-G-G-AGT-GTCC-AA-AGCTC-T-T-TGT-T-TGA-T-GGCAT-CT-CTG--T----
T-TA-CA--GAGTTTACACT-T-T-AATA-T-CA-A-CCT-GT-TTCCTCCTCC-T--CCTT-C-TC-CTC----CT-CCT-CC-GT-GAC-CT-CC---TC---
-CTAC-C-TCT-T-T-CT-C-CT-GAGA-AA-CT-T-C--GC-CC--CAG--CGGTGCGGA-GCGCCCTGC-GCA-GC-CGG-GGAGG-T-A-CGCA--C----
C-CGC-C-----G-CGC-----CA------CA-----
clustal w Result : 1328.0
-TTACAGAGTTTACA-----CTTTAATATCAACCTGTTTCCTCCTCCTCCT
   -----CGCACTGCAGGGGCGCCAGATTTGGCGGGAGGGG
GCGACGACCCGGCCAGCC----CGGCACCCGCGGGCGGCAGCC----AGGGCGACGCGGA
     -GCAGGCAGGCGGC-GGGCAGCGGGAGGCGGCAGCCCGGTCGGTCCCCGCGGC
-CTACCTCTTTCTCCTGAGAAACTTCGCCCCAGCGGTGCGGAGCGCCCTGCGCAGCCGGG
   -CTCCCTGGCCGGCTCCCTGCCGGGCCTCCCGGGCCTGGTGCTAGGGCACCGCGG
   -----TGGGCGCGGGTCGGCGGCCGCGACCCGGGAGCGGGTTTGCTCAGGAAAAG
  ------
```

```
SAGA Result : 863.5
C-C-C-AC-T-GG-C-A-AGCT-T-G-AGGAGCCAGG-CT-GCCAGTCGGGAG-AT-TCGG------CC-C-AGTGT-T-C-C---CA--CTGGA-GAGGGC-
GGCA-A--G-TG-CCCGGGCGAT-CACCTCG--CCT--GCG-T-T-CG-GGAGA-TAT-A-C-C-TC-CGC-C-C-C-C-C-C-CCGC-CA-GGAG-GG---
T-GA-AAAGATGGCCCC-A-G-GAG-CCAG-C-CGG-CTGG-GA--CA-AGGC-GGAGTGAG-AG-G-A-CA-GG-CT-G-G-G-----GCC-GGGGGC-----
-GCT-G-GGCTGTCCCGGGC-AG-C-CCTCC-TCCG-GGCA-A-G-CC-GGAGC--AG-G-G-G-G-GA-T-T-G-G-GA-G--CGCTCG-GG-GCGG---
GCC-CGCG-GTGGCCC-GGGGCGGTGG-CGC-CC-GGC-C-GGA-GAG-GGTG-G-G-GCG-G-AGA-CGC-C--GC-C-TG-TACT--TC--C-C----
-CT-T-C-G-CCG-C-T-A-GCT-C-T-A-CAACA-GCCTGAT-T-T-CC-C-C-GA-AATG-A-CGGC-ACGC-AGC-C--GGCCAA-TGG-G-CGC--CGC
GCG--GCT-GT--CCG-GGGGCGG-GGCCGG-CCAGGG-CTGGG-GAA---TC-C-C-GCT-A-AGT--GT-T-TGG-A-T-TGCT-CG-GTG-GCGC-----
CG--C-----GCCCT-G-GC-----
clustal w Result : 1179.0
-----CTTCGCCGCTAGCTCTACAACAGCCTGATTTCCCCCGAAATGACG
---CCCAC----TGGCAAGCTTGAGGAGCCA-GGCTGCCAGTCGGGAGATTCGGC
---GCGGCTGTCCGG------GGGCGGGG-CC-GGCCAGGGCTGGGGAATCCC-GC
-GCTGGGCTGTCCCGGGCAGCCCTCCTCCGGGC-AA-GCCGGAGCAGGGGTGGATTGGGA
TGAAAAGATGGCCCCAGGAGCCAGCCGGCTGGG-AC-AAGGCGGAGTGAGAGGACAGGCT
GCCCGCGGTGGCCCG--GGGCGGTGGCGCCCGG-CC-GGAGAGGGTGGGGCGGAGACGCC
     ----GGCAA--GTGCCCGGGCGATCAC-CT-CGCCTGCGTTCGGGAGATATACC
-----C--GC------
```

CACCCCTCGCCCGCACCCCTGGCCCAAAACAACTGGCCAGGTTCCCTGGC CTCCCGGGTCCCTGCATCCCCGCATCCCGTCCGCAGCCGTGAACTTGA GCCCCCCTCCATCAGAGGTTGCGAGCGTCCGCCCCGCTCGCGGCAGCCACC GTCACTAGACAGTCAAACCCCCAAGACGTCAGCCCACAATGCACCGGGCGG GCCGGGAAAAACGGCCCGGGGAGGGGACCGGGGAAGAGAGGGCCGAGAGG CGTGCGGCAGGGGGGGGGTAGGAGAAAGAAGGCCCCGACTGTAGGAGGG CAGCGGAGCATTACCTCATCCCGTGAGCCTCCGCGGGCCCAGAGAAGAAT CGACGCGAGCCAATGGGAAGGCCTTGGGGTGACATCATGGGCTATTTTTA GGGGTTGACTGGTAGCAGATAAGTGTTGAGCTCGGGCTGGATAAGGGCTC

```
SAGA Result : 1013.5
C-A-CC-CCTC-GC-CC-G-CAC-C-CCTGGCCCAAAACA-A-CT-GG-CCAG-G-T-T-C-CCT-GGC----
C-TCCCGGGT-C-C-C-T-GCA-T-CCCC-C-GC-AT-C-C-CCGT-CCGC-A-GC-C-GT-GAA---CT-TGA--
GCC-CCCCTC-CA-T--CAG-A-G-G-TT-GCG--AGC-G-TCC-G-CC-C-GCTCG-CGGC-A-G-CCA-C-C-
G-T-CA--CTA-GA-CA-GTCAA-A--C---CCCAAGACG-T-CA-GC-CCAC-A-A-TGCACCG-GGCGG-----
CA-G-CG-GAGC-A-TT--AC--CT-CATC-CCGTGA-GC--CT---CCGCGG-G-CCCA-G-AGAAGAA-T----
C-TTCTAGG---GTG-G-AG-T-CTCCAT-GGTGACGG-GCG-GGCCCG-CC-C---CCCT-GAGAG--
CGA-CGCGAGCCAA-T-GGGA-AGGC-CTT-GGG-GTG-AC-ATCATGGGCT---ATTT-TT-----A-----
G----GGGTT-G-ACTGG--TA-G-CA--G-ATA-AG-T-GTT--G-AGC-TCGGGCTGGATAAG-G-GC-TC--
clustal w Result : 2229.5
   ----CTTCTAGGG-T-GGAGTCTCCATGGTGACGGGCGGGCCCGCCCCCCTGAGAG-
GCCGGGAAAAACGGCC-C-GGGGAGGGGACCGGGGAAGAGAGGGCCGAGAGG------
     CGTGCGGCA-G-GGGGGGGGGGTAGGAGAAGAAGGGCCCGACTGTAGGAGGG-
---GGGGTTGACTGGTAGCAGATAAGTGTTGAGCTCGGGCTGGATAAGGGCTC-----
    --CGACGCGA----GCCAATGGGAAGGCCTTGG-----GGTGACATCATGGG
  ----CACCCCTCGC-CCGCACCCCTGGCCCAAAACAACTGGCCAGGTTCCCTGGC---
    --GCCCCCCTC-CATCAGAGGTTGCGAGCGTCCGCCCGCTCGCGGCAGCCACC--
-----CTCCCGGGT---CCCTGCATCCCCGCATCCCCGTCC--GCAGCCGTGAACTT
      -GTCACTAGA-CAGTCAAACCCCAAGACGTCAGCCCAC--AATGCACCGGGCGG
    -----CAGCG-GAGCATTACCTCATCCCGTGAGCCTCC--GCGGGCCCAGAGAA
```

AGAAACCCTGTCTCGAAAAACAAACAAACAAACCTAAAAATAATAATGAA CAAATAAACCCTCCGTCTACCAGATTAGAGGAGGGAGTGGGAAGTGCTCT AACCTGGAAGGAGAACGCTGCTTATAGTGTAGACCCTTTTCCTTTAAGAC CAAAGACCGACTCAAATTTCCCAGAAACTGTAGAGGTAGCTCCAAGCTCC ACCCTTCTAGTCTTCAGGTTTTGTTGTCGCGCCGTAGCGTTTAAATTTCG CGCGCTCTACACGGATTGGCCATTTCCAGGGGTGAGCAGCCTCCGGCTTG CGGCCACTCCGCCTTGTCTTCTTCCCCGGGTAAAAGGTTTTCAAATTGGA CCAATCACTGGGCGCGTTCTCTCAGCTAAGGCGCGTCACCGAAGGGTTAA TTGCAACCAACCAGAGGTGGGTATTAAAAAAGGAACCAATCAGGAGGAAG

```
SAGA Result : 870.5
A--GA-AACC--CT--GT-C--TC-GA-A-AA-ACAAACAAAC-A-AACCTAAAAATA-AT--A-ATGAA--
CAAAT-AAACC--CTCC-G--TC-TA-CCAGAT-TAG-A-GG-AG-GG-AG-TG-G-GAAGTGCTCT-----
AA-CCTG-G-AAGG-A-G-A--AC-GC-TGCTTA-TA-----GTGTA-GAC-CCT-TTT-CCTTTAAGAC
CAAAG---ACCGACTCA-A-AT-----T-TC-CCAGAAAC-T-GT-AGAGGTAGCTC-C--AAGCTCC---
AC-C-CTTCTAGTCT-TCAGGT-TTT-G---TTGTCGC-GC-C-GTAGC-GT----T-TA-A-A-TTTCG--
CGCG-C-TCTACA-C-GGA-T-T-GG-C-CATTTC--C-AG-GGGT-GA-GCAGCCTC-CGGCTTG-----
CGGC-C-A-C--TCC-GC-C-T-TGTCTTC-TT-C-C--CCGGG-TA-AAAGGT-T-T-TCAAA-TTGGA-
C---CAAT-C-ACT--G-GGC-GC-GTT-CTCTC-AG-C-TAAGGC-G-CG---TCAC-CGAAGGGTTAA--
T-TGCAAC-CAACC-AGAGGT-GG-GTA--T-TA-AA-A--AAGGA-A-C--CAA-TC-AG-GAGG--AAG-
CGGG-G-G-CGTGTCCAGG-G-A--GT-TTA-TA-A-G-GCCGGGCTT-----GGCC-G-G-CC-AGGCGA-
clustal w Result : 1934.0
  -----AC--CCTTCTAGTCTTCAGGTTTTG-TTGTCGCGCCGTAGCGTTTAAATTTC
  -----A--ACCTGGAAGGAGAACGCTGCT--TATAGTGTAGACCCTTTTCCTTTAA
--CAAATAAA--CCCTCCGTCTACCAGATTAGA-GGAGGGAGTGGGAAGTGCTCT----
--CAAAGACC--GACTCAAATTTCCCAGAAACT-GTAGAGGTAGCTCCAAGCTCC----
----CCAAT--CACTGGGCGCGTTCTCTCAG--CTAAGGCGCGTCACCGAAGGGTT-AA
-CGCGCTCTA--CACGGATTGGCCATTTCCAGGGGTGAGCAGCCTCCGGCTTG-
---CGGCCAC--TCCGCCTTGTCTTCTTCCCCGGGTAAAAGGTTTTCAAATTGGA----
----AGAAA--CCCTGTCTCGAAAAACAAAC-AAA-CAAACCTAAAAATAATAATGAA
TTGCAACCAA--CCAGAGGTGGGTATTAA-----AAAAGGAACCAATCAGGAGGAA
```