

# Поравнање више секвенци

# Циљеви часа

- Разумевање ситуација у којима је поравнање више секвенци корисно
- Како да се исправно направи поравнање више секвенци
- Разумевање предности и мана прогресивног алгоритма за поравнање више секвенци

## План часа

- Зашто се прави поравнање више секвенци?
- Како се бирају одговарајуће секвенце
- Како ради прогресивни алгоритам за поравнање
- Интерпретација поравнања више секвенци

# Шта је поравнање више секвенци ?

- Поравнање више секвенци
- MSAs = multiple-sequence alignments
- Постоје два циља MSA:
  - Утврђивање који делови секвенце одговарају једни другима
  - Налажење позиција које су очуване
- Кључни кораци у MSA су:
  - Исправно одабирање секвенци које се поравнавају
  - Коришћење одговарајућег MSA метода алгоритма
  - Интерпретација поравнања

# Кратко о очуваним позицијама

- Аминокиселине случајно мутирају
- На мутације онда делује позитивна селекција (бивају прихваћене) или негативна селекција (бивају одбачене)
- Ако је мутација штетна, на њу делује негативна селекција
  - Нестаје из генома
  - Никада је не видите
- Мутације битних позиција (нпр каталитички активна места) су готово увек штетна
- Можете да препознате битне позиције пошто готово никад не мутирају!
- MSAs открива *очуване* позиције

# Пример очуване позиције:

```
CLPP_ECOLI E.col (40) ERVIFLTGQV---EDHMANLIVAQMLFLEAENPEKDIYLYINSPGGVITAGMSIYDTMQFIKPD---VSTIC (105)
CLP1_MYXXA M.xan(26) DRIIMLGTFV---NDDVANIIVAQLLFESEDPDKGINLYINSPGGSVTAGLAIYDTMQYVKCP---VSTIC ( 91)
21228980 M.maz (27) MISLFGLPAYQSIDEEDAQVLRWIRKY----RDYPLELILHTPGGQLHASIQIARALKNHHPKK---TRVLI ( 92)
15643678 T.mar (58) SISFLGFPVRRYIDIEDSEEILRAIKLTP----SDMPIDLILHTPGGLVLAAEQIARALKMHKGGK---VTVFV (123)
15668307 M.jan(64) SIGLFGIPVYKFTIEDSEEILRAIRAAP----KDKPIDLIHTPGGLVLAATQIAKALKAHPAE---TRVIV (124)
18976612 P.fur (59) SIGFFGIPVYKFTISIEDSEEVLRAIRMAP----KDKPIDLIHTPGGLVLAATQIAKALKDHPAE---TRVIV (124)
22972030 C.aur (53) TMSLLGFPVLRVINIEDSEAVLRAIKMTD----RDIPIDLILHTPGGLVLAAEQIARALKKHAAK---VTVFV (118)
23050732 M.bar (75) AISLFGIPAYQYIDEEDAQILRWIRKY----KDYPLELILHTPGGQLHSSIQIARALRRHSN---TVVIV (159)
15964138 S.mel (50) HVARVAVTGLIQ---DDRELVERLERIADN---QSVKALIVTISSPGGTTYGGEVYKLRKVAEKKE---VVSVDV (116)
17934547 A.tum(27) AIMAGGNQFRPALNLASYAPLLEKAFVAKDA---PAVAISLNSPGGSPVQAPNINRIROLAERKDKKVLIFV ( 96)

CLPP_ECOLI E.col (106) MGQAASMGAFLLTAGAKGKRFCLPNSRVMIHQEGGY-----QGQATDI----- (147)
CLP1_MYXXA M.xan( 92) VGQAASMGALLLLGAGKGRYALPNSRPHIQPLGGA-----QGQATDI----- (133)
21228980 M.maz ( 93) PHYSMGGTIIALAADI-IVNDRDAVIGPID-PQVGDPIRGVFPMSNIHAAETKK-EDADDSTLVMS----- (156)
15643678 T.mar (124) PHYAMSGTIIALAADI-IINDENAVLGPVD-PQIGN-----YPAPSILAAVKKKDVNEVDDQTLILA----- (184)
15668307 M.jan(130) PHYAMSGTIIALAADI-IINDENAVLGPVD-PQLGQ-----YPAPSIVKAVEQKQKADKADDQTLILA----- (190)
18976612 P.fur (125) PHYAMSGTIIALAADI-IINDPHAVLGPVD-PQLGQ-----YPAPSIKAVEQKQKAEKVDDQTLILA----- (185)
22972030 C.aur (119) PHYAMSGTIIALAADI-IVNDRDAVIGPID-PQLGQ-----HPAASILSVLERKPLSEIDDETMMMA----- (179)
23050732 M.bar (140) PHYSMGGTIIALAANE-IVNDRDAVIGPID-PQIGDFIRGMYPAPSWIYAAETKK-EKADDTLVMS----- (204)
15964138 S.mel (117) RTLAASAGYLIALAGDR-IVAGETSITGSIG-VIFQY-----PQVKTLMDKLGVSLSEIKSEELKAEPSPFHPSPS (184)
17934547 A.tum( 97) EDVAASGGYMIALAGDE-IIADPTSIVGSIG-VVSGG-----FGFPEMLRKI GVERRVYTS ENKVVLDPPFQPEK (164)

CLPP_ECOLI E.col (148) ----EIHAREILKVKGRMNEIMLALHT-----GQSLEQIERDT-----ERD-RFLSAPEAVEY (196)
CLP1_MYXXA M.xan(134) ----DIQAKEILRLRSYINGLIVKHT-----GHTIERIEKDT-----ERD-YFSAEDARQY (182)
21228980 M.maz (157) ----DISRKALRLTRNVAKELLEGGKIQPD-GKEDRLEEVVEKLVSG-EMINSTPLSAREAKEL (213)
15643678 T.mar (185) ----DIAEKAI RQVKEFVVEILSDKV-----SKEKAEKIADKLCSG-YWTHDYPLNVEKLREM (237)
15668307 M.jan(191) ----DIAKKAINQVQNFVYNLLKDKY-----GEEKAKELSKILTEG-RWTHDYPIITVEEAKEL (243)
18976612 P.fur (186) ----DVAKKAIKQVQDFLYDLLKDKY-----GEEKARELAQILTEG-RWTHDYPIITVEHAREL (238)
22972030 C.aur (180) ----DIAEKAI RQVKRTVCELLRDKM-----PVERAEVVAHTLASG-VWTHDYPIITVSEAREL (232)
23050732 M.bar (205) ----DVSRRKALKFTRNVAKELLEGGKIQGPPAGESRLDEVVEKLVSG-EMINSTPLSAGEAKKI (262)
15964138 S.mel (185) DEARAMIQAMIDDSYGNFVDLVAERRK-----LPRPEALALADGRI FTGRQALEGKLVDEL (240)
17934547 A.tum(165) EGDIDYLSLQVEIHNVFIDMVKMRRG-----SKLK--GDDALFSGLFWTGMRGLDLGLIDGL (220)
```

Активно место



# Зашто правимо MSA ?

- Главни разлог је да би га користили у даљим применама
- MSA се често користи у прављењу различитих модела
- Следећа 2 слајда дају најчешће примене MSA

# Примене MSA

Примена	Процедура
Екстраполација	Одређивање функције протеина
Филогенетска анализа	Прављење филогенетског дрвета
Идентификација образаца (pattern)	Открите битне делове секвенце
Идентификација домена	Претворите поравнање у профил домена
Идентификација регулаторних елемента у ДНК	Анализа промотера
Предвиђање структуре	Предвидите секундарну структуру РНК и протеина
PCR анализа	Предвиђање PCR прајмера



# Одабирање одговарајућих секвенци

- Први корак у MSA је да одаберете одговарајуће секвенце
- Два главна фактора у одабирању секвенци:
  - Број секвенци
  - Природа секвенци
- Разуман број секвенци: 15 до 50
  - Идеалан за већину метода
  - Мања поравнања се лако приказују и анализирају
- Типови секвенци
  - Добро селектоване секвенце  $\Leftrightarrow$  информативно поравнање

# Нека упутства за одабирање одговарајућих секвенци

<i>Problem</i>	<i>Diagnostic</i>
<b>Proteins or DNA</b>	Use proteins whenever possible.
<b>Many sequences</b>	Start with 10~15 sequences, 50 at max
<b>Very different sequences</b>	Avoid sequences very different from the rest of the set
<b>Identical sequences</b>	Avoid identical sequences; they never help
<b>Partial sequences</b>	Can trigger errors. Avoid them
<b>Repeated domains</b>	Can trigger errors. Extract the domains with Dotlet and align each domain.

# Нека упутства за одабирање одговарајућих секвенци

Проблем	Решење
Протеини или ДНК	Користите протеине кад год је могуће
Колико секвенци	Почните са 10-15, максимално 50
Веома различите секвенце	Избегавајте секвенце веома различите од осталих
Идентичне (јакко сличне) секвенце	Никад не помажу
Парцијалне секвненце	Избегавајте их
Домени који се понаваљају	Извуците домене који се понаваљају помоћу dotlet-a

# Одабирање одговарајућих секвенци

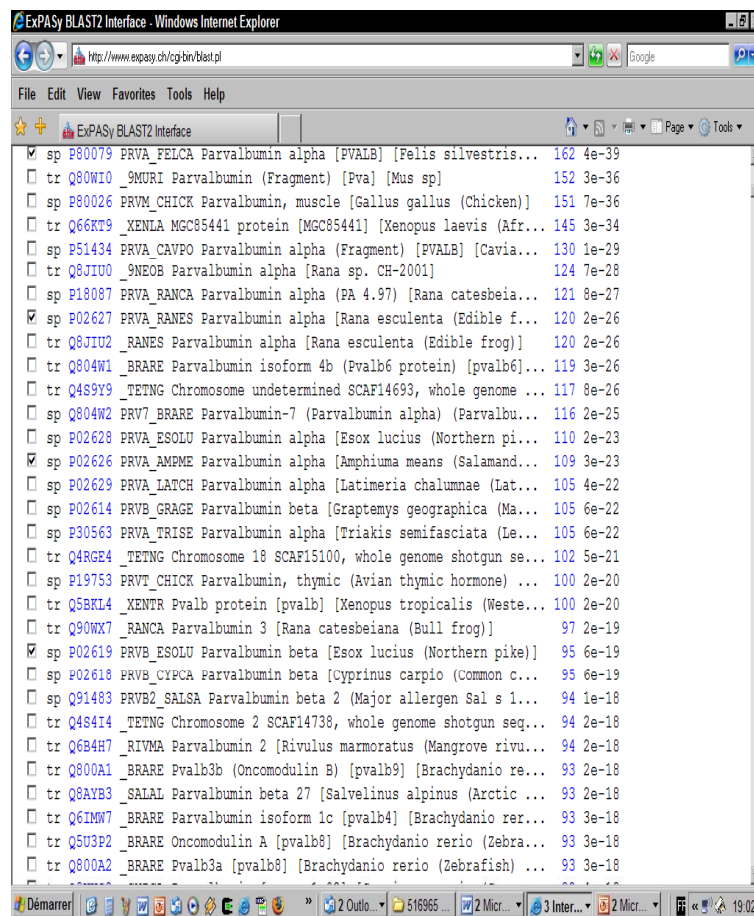
- Поравнање је корисно ако . . .
  - Су секвенце које поравнавате исправно изабране
- Да би добили резултате, секвенце морају да буду
  - Не превише сличне
  - Не превише различите
- Секвенце које су веома сличне . . .
  - Су једноставне да се добро поравнају
  - Нису превише информативне  $\Leftrightarrow$  лоша предвиђања
- Секвенце које су веома различите . . .
  - Су тешке да се поравнају
  - Су веома информативне  $\Leftrightarrow$  потенцијално добра предвиђања

# Сакупљање секвенци помоћу BLAST- у

- Најпогоднији начин да селектујете секвенце је да користите BLAST сервер
- Неки од BLAST сервера су интегрисани са методама за поравнање више секвенци:
  - [www.expasy.ch](http://www.expasy.ch)
  - [srs.ebi.ac.uk](http://srs.ebi.ac.uk)
  - [npsa-pbil.ibcp.fr](http://npsa-pbil.ibcp.fr)

# Сакупљање секвенци помоћу BLAST-а

- Равномерно селектујте секвенце од врха до дна
- Идеја је да имате равномеран прелаз између секвенци

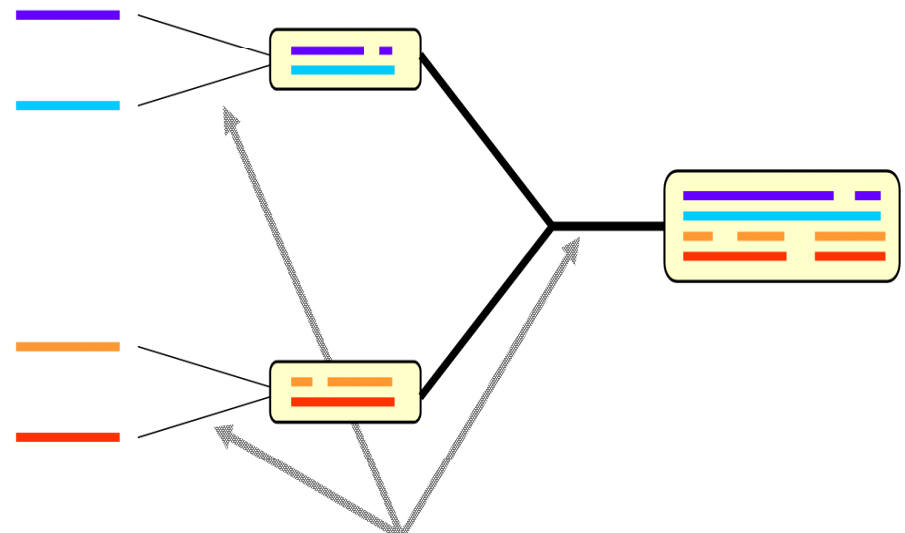


# Поравнање секвенци

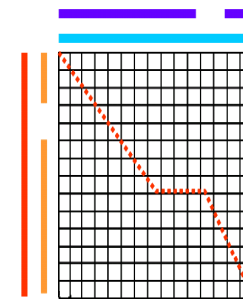
- Генерално веома тежак проблем
- Сви методи су апроксимативни
- Методи за поравнање користе прогресивни алгоритам
  - Секвенце се упоређују две по две
  - Прави се дрво које дефинише процес упоређивања
  - Секвенце се поравнавају по реду које утврђује то дрво

# Прогресивни алгоритам

- Секвенце се групишу по сличности (прави де дрво)
- Секвенце се поравнавају 2 по 2
- Добијена поравнања се онда даље поравнавају 2 по 2



Dynamic programming using a substitution matrix





# Прогресивни алгоритам (наставак)

- Главна предност му је брзина
- Главна мана му је похлепа
  - Не враћа се на секвенце које је поравнао на почетку
  - Ране грешке се не отклањају
- Сакупите податке са доста прелазних секвенци
- Замислите да је свака секвенца део каменог моста преко реке:
  - Није битно колико је широка река, ако је камење довољно близу
  - Није битно колико су секвенце међусобно различите, докле год свака има блиског рођака



# Селекција метода

- Постоје многе MSA методе
- Готово сви користе **прогресивни алгоритам**
- Сви методи су апроксимативни
- Ниједан није гарантован да води ка најбољем поравнању
- Сви постојећи методи имају добре и лоше стране
  - **ClustalΩ (наследник ClustalW)** је најпопуларнији (21,000 цитата)
  - T-Coffee и ProbCons су тачнији али спорији
  - MUSCLE је веома брз, идеалан за веома велике базе

## Неки коментари о MSA . . .

- У пракси, прављење MSA није егзакна наука
- Одабир секвенци је најбитнији корак
  - Ако поравнање изгледа лоше, промените секвенце
- У пракси, прављење доброг MSA је итеративан (try-again) процес
  - Селектујете пар секвенци
  - Поравнате их
  - Визуелно оцените квалитет MSA
  - Додате или одузмете секвенце
  - Поново их поравнате

## Интерпретација MSA

- Користите обрасце конзервације ClustalW
  - ‘\*’ Комплетно конзервирана позиција
  - ‘.’ Високо конзервирана позиција
  - ‘:’ Конзервирана позиција
- При интерпретацији вашег поравнања, увек користите биолошку информацију (ако је на располагању)

# 1. Направите поравнање протеинских секвенци сигма фактора из ECF групе (сигма70 фамилија), датих у фајлу ECF\_sigma\_1.txt

Напомена: ECF сигма фактори учествују у регулацији експресије бактеријских гена под стресним условима, при чему за интеракцију са ДНК промотора користе домене  $\sigma^2$  и  $\sigma^4$ .

Упутство: Користите програм Clustal $\Omega$  (default параметри) и на поравнању уочите границе домена  $\sigma^2$  и  $\sigma^4$ .

Додатак: Који од поравнатих региона одговара  $\sigma^2$ , а који  $\sigma^4$  домену? Да би ово са сигурношћу утврдили, изаберите неку од секвенци (нпр. прву) и пустите на њој CD search. Искористите резултат CD search-а да кратко прочитате о функцији ова два домена.

# Локално поравнање више секвенци

- Већина MSA програма подразумева да су ваше секвенце повезане целом дужином
- Када ова претпоставка није тачна, прогресивни приступ не ради
- Једина алтернатива је да урадите локално поравнање више секвенци

# Методи за локално MSA

- Gibbs Sampler
  - Игнорише неповезане делове ваших секвенци
  - Идеалан за налажење мотива у ДНК
- Методи за локално MSA
  - Тражи кратке мотиве који су очувани у секвенцама
  - Секвенце не треба да се глобално поравнају
- The most popular motif-discovery methods:
  - MEME, Bioprosector

3. Направите вишеструко локално поравнање ДНК секвенци, датих у фајлу MLSA.txt, при чему је дужина траженог мотива 5bp

Упутство: Користите програм MEME, <http://meme-suite.org/tools/meme>, MLSA.txt унесите под **Input the primary sequences**. Користите default параметре, осим:

- 1) Под **Select the site distribution**, користите **oops** опцију.
- 2) Под **Input job details** унесите вашу мејл дресу, ово је опционо, али типично ћете испробавати различите дужине мотива и њихове дистрибуције у секвенцама (види тачке 4 и 1 изнад), па су резултати овако прегледнији.
- 3) Под **Select number of motifs** одаберите **1**
- 4) Под **Advanced options** ставите и **minimum width** и **maximum width** на 5; ова опција одређује дужину мотива коју тражите у секвенцама – у овом случају тражите да мотив буде 5bp дугачак.



# Задаци вежба

2. Исто као задатак означен са 1 (глобално поравнање више протеинских секвенци помоћу ClustalΩ), само сада поравнајте протеинске секвенце сигма фактора датих у фајлу ECF\_sigma\_2.txt. Секвенце такође припадају ECF групи сигма фактора, само су другачије селектоване. Упоредите поравнање са резултатима добијеним у задатку 1 (информативност поравнања, тј. могућност уочавања домена  $\sigma^2$  и  $\sigma^4$ ). Коментаришите у којој мери селекција секвенци може да утиче на поравнање, односно повежите са дискусијом у предавању.

4) Исто као задатак означен са 3 (локално поравнање више секвенци помоћу MEME-а), само сада урадите поравнање и са дужином мотива 4 и 6. Колико је мотив који налазите робустан у односу на овакву промену дужине? Такође, тестирајте и робусност у односу на промену у **Select the site distribution**, нпр. уместо **oops** одаберите **zoops** опцију. Да ли је добијени мотив и даље робустан? Приметите да logo (интуитивну визуелну репрезентацију) мотива можете да видите под **Discovered motifs** (кликом на плаву стрелицу усмерену надоле, лого можете да увећате) – помоћу логоа вам је најлакше да процените да ли добијате исти мотив (скраћен, померен за 1bp налево или надесно итд.) Да ли можете да препознате на шта се односи нађени мотив? У вези са тим, ако су секвенце које претажујете интергенски региони узводно од бактеријских гена (прецизније узводно од оперона), зашто очекујете да се овај мотив појави у свакој од претраживаних секвенци (односно у готово свакој)?