

آرثر م. ليسك

مقدمة في

البيومعلوماتية

عرض: د. / محمد باسم عاشور

إجتهادات حديثة حول العلم والمستقبل
« عروض »
السلسلة

ISO
9002



المكتبة الأكاديمية

شركة مساهمة مصرية - القاهرة

EBSCO Publishing : eBook Arabic Collection Trial - printed on 4/9/2020 6:30 AM via MINISTÈRE DE L'ÉDUCATION NATIONALE, DE LA

FORMATION PROFESSIONNELLE

AN: 844665 ; Lesk, Arthur M., . ; : Introduction to bioinformatics

Account: ns063387

مقدمة في البيومعلوماتية

Introduction to Bioinformatics

تأليف

Arther M. Lesk

الناشر

Oxford Univ. Press, 2002

عرض

د. محمد باسم عاشور

أستاذ علم السمية - جامعة الزقازيق

الناشر



المكتبة الأكاديمية

ش.م.م

٢٠٠٦

هذه الكراسة تقدم عرضاً تفصيلياً لكتاب :

Introduction to Bioinformatics

Arther M. Lesk

Oxford Univ. Press, 2002

حقوق النشر

الطبعة الأولى ٢٠٠٦م - ١٤٢٥هـ

حقوق الطبع والنشر © جميع الحقوق محفوظة للناشر :

المكتبة الأكاديمية

شركة مساهمة مصرية

راس المال المصدر والدفوع ١٨,٢٨٥,٠٠٠ جنيه مصرى

١٢١ شارع التحرير - الدقى - الجيزة

القاهرة - جمهورية مصر العربية

تليفون : ٧٤٨٥٢٨٢ - ٣٣٦٨٢٨٨ (٢٠٢)

فاكس : ٧٤٩١٨٩٠ (٢٠٢)

لا يجوز استنساخ أى جزء من هذا الكتاب بأى طريقة
كانت إلا بعد الحصول على تصريح كتابى من الناشر .

هي الثالثة في مشروع "الكراسات"، الذي تصدره "المكتبة الأكاديمية". والكراسات تعنى بمحورين كبيرين: العلم والمستقبل. لذلك فقد حملت السلسلة الأولى عنوان "كراسات مستقبلية"، وقد بدأ ظهورها عام ١٩٩٧. وفي عام ١٩٩٨ ظهرت السلسلة الثانية تحت اسم "كراسات علمية". وقد فكرنا في البداية أن تضم السلسلتين، بجانب التأليف والترجمة، عروضاً مطولة لبعض الإصدارات المهمة، التي لا تلاحقها حركة الترجمة. إلا أن أنشط أعضاء أسرة الكراسات، وللكراسات أسرة ممتدة ترحب دائماً بالأعضاء الجدد، أقول أن أنشط الأعضاء الصديق الدكتور محمد رؤوف حامد، الأستاذ بهيئة الرقابة الدوائية، اقترح أن تصدر العروض في سلسلة خاصة بها. وقد كان اقتراحاً موفقاً كما أرجو أن يوافقني القارئ.

والكتب المختارة للعروض في السلسلة لاتأتى فقط من اقتراحات هيئة التحرير، حيث قدم أعضاء الأسرة مقترحاتهم التي حظيت بالترحيب، والباب مفتوح لكل من يرغب في المشاركة. وإذا كانت السلسلة قد بدأت بمجموعة من الكتب الصادرة بالإنجليزية، فإننا نطمح أن تشمل العروض القادمة كتباً تصدر في لغات أخرى، لانتشملها عادة خطط الترجمة كاليابانية والروسية والصينية، بالإضافة إلى الفرنسية والألمانية. فرغم أن الأخيرتين أكثر حظاً نسبياً، إلا أن كم المترجم والمعروض لا يقارن بما يتم بالنسبة للإنجليزية.

والحديث عن "العروض" يذكرنا بالجهود السابقة، التي لانكرها، بل نحاول أن نكمل مسيرتها. فبالنسبة للعروض الموسعة، تذكر جهود الهيئة العامة للإستعلامات بالنسبة للمجالات التي تهتمها. كما أن العروض المتوسطة، التي أصدرتها هيئة الكتاب في التسعينيات، ضمن سلسلة "تراث الإنسانية" لا يمكن إغفالها. وهما مثلات يقصد بهما الإعراف بفضل سبق، دون أن ندعى الحصر. وإن كنا، في الوقت نفسه، نظن أن السلسلة الحالية هي الأولى التي تعنى بالعروض التفصيلي للكتب.

هذه الكراسة

هي الأولى التي يمكن للقارئ المتخصص أن يطالعها وهو متصل بالإنترنت ليمارس بعض تطبيقات البيومعلوماتية الموضحة بها. لكنها أيضاً موجهة للقارئ العام، الذي يريد أن يلم بهذا المجال الجديد، الذي ذكر أنه قد حول البيولوجيا (علم الحياة) إلى علم معلوماتي. ولأن المجال يجمع بين البيولوجيا والمعلوماتية فالبيولوجي المهتم بتطبيقاته يحتاج إلى الإلمام بالقاعدة المعرفية والمهارية المعلوماتية التي تمكن من ذلك. والمشتغل بالمعلومات يحتاج بدوره إلى مزيد من المعرفة عن البيولوجيا التي سيوظف تخصصه فيها. والقارئ العام، الذي يسعى إلى تثقيف نفسه علمياً، يهمله أن يفهم كيف تساهم البيومعلوماتية في تشكيل صورة المستقبل في الطب والزراعة والبيئة... إلخ. لذلك رحبت السلسلة بالعرض المتميز لكتاب أرثر ليسك "مقدمة في البيومعلوماتية"، الذي أعده الدكتور محمد باسم عاشور، أستاذ علم السمية ووكيل كلية الزراعة جامعة الزقازيق، الذي سبق أن قدمنا له كراسة مستقبلية عن التكنولوجيا الحيوية الزراعية.

أ. د. أحمد شوقي

يناير ٢٠٠٦

البيومعلوماتية مجال علمي تطبيقي جديد، زادت أهميته والحاجة إليه مع توفر الكم الهائل من المعلومات والتي نتجت عن مشروع الجينوم البشري، وغيره من جينوم كائنات عديدة، وكذلك ضرورة تشخيص فيروسات تسبب أمراض خطيرة في الكائنات وما يتطلبه من تحديد تتابعات والتعرف على تراكيب، وإنتاج أمصال، وتصميم دواء، وحلول لمشاكل المقاومة للمضادات الحيوية والدواء والمبيدات.

تقدم البيومعلوماتية مجموعة من المهارات المتعلقة بأدوات جمع المعلومات والتنقيب عن البيانات وبناء المعرفة في البحوث الجارية والتطبيقات الصناعية والأكلينيكية وغيرها. وذلك باستخدام برامج كمبيوتر لعمل استنتاجات من أرشيف بيانات البيولوجيا الجزيئية الحديثة، والربط بينها والوصول الى تنبؤات مفيدة. ويتحقق ذلك ليس فقط عن طريق تطوير وزيادة السعة التخزينية ولكن أيضا بابتكار النماذج الحسابية والصور التقديرية لمعالجة البيانات.

ومجال البيومعلوماتية هو نتاج التزاوج بين علوم البيولوجيا وعلوم الكمبيوتر وهندستها ويحظى باهتمام المشتغلين بتلك العلوم. لذلك من المهم بالنسبة لعلماء البيولوجيا دراسة الكمبيوتر وكذلك تدريس البيولوجيا للمتخصصين في علوم الكمبيوتر. حيث يتطلب هذا المجال تكامل وفهم للخلفية البيولوجية مع اكتساب مهارات ضرورية في مجال الكمبيوتر.

ويهدف كتاب "مقدمة في البيومعلوماتية الى بناء وتطوير خلفية عن مجال البيومعلوماتية بدون الاعتماد على طرق معقدة في علوم الكمبيوتر أو مهارات برمجة متقدمة. وذلك بدعم وتشجيع استخدام الأدوات الحسابية والصور التقديرية بطريقة مرضية وشيقة. وكذلك يتيح الكتاب للقارئ التفاعل مع موضوع الدراسة، وتطوير المادة العلمية طبقا للحاجة مع توفير مصادر لمواقع الكترونية متنوعة.

دكتور محمد باسم عاشور

الفصل الأول

مقدمة

علم البيولوجى من العلوم التى تعتمد على المشاهدة وحديثاً أصبح من العلوم الاستنتاجية حيث تغيرت طبيعة النتائج. ومن المعروف للجميع أن جميع المشاهدات البيولوجية الأساسية تشتمل على سرد لروايات بدرجات متفاوتة من الدقة. إلا أنه فى الفترة الأخيرة أصبحت البيانات ليست فقط أكثر كما ودقة، ولكن متحفظة وحكيمة كما هو الحال فى بيانات تتابعات النيكلوتيدات فى الأحماض النووية.

ومن الممكن الآن تقدير تتابع الجينوم لكائن منفرد أو نسخة منه ليس فقط كاملاً بل على نحو صحيح ودقيق. ومع أنه ليس من الممكن تجنب الخطأ التجريبي إلا أنه منخفض للغاية بالنسبة للتتابع الجينومى.

وكذلك من الخصائص الواضحة لبيانات البيومعلوماتية أنها ذات كم هائل جداً. فمثلاً يحتوى بنك البيانات لتتابع النيكلوتيد على 16×10^9 من القواعد (16 Gbp). لو استخدم الحجم التقريبي للجينوم البشرى (3,2 * 10⁹ حرفاً) كوحدة فإن ذلك يعادل (2 bugs, an apt name). وللمقارنة المرجعية فإن 1 bugه يقابل عدد الحروف التى ظهرت فى اعداد نيويورك تايمز لمدة 6 سنوات. تحتوى قاعدة البيانات لتراكيب الجزيئات الكبيرة (Macro molecules) (الأبعاد الثلاثية للبروتينات بمتوسط طول 400 حمض أمينى) على 16000 قيد.

شجع الكم الهائل المتاح من البيانات العلماء الى السعى نحو أهداف قياسية وطموحة مثل:

- التأكيد على قدرة العلماء على الرؤية الواضحة والكاملة للحياة. بمعنى الفهم المتكامل لبيولوجيا الكائنات كأنظمة مترابطة ومعقدة.

- ايجاد وربط علاقة التتابع والتركييب ثلاثى الأبعاد والتداخلات والوظيفة لبروتينات وأحماض نووية منفردة وكذلك لمعقد البروتين - الحامض النووى.
- استخدام بيانات لكائنات معاصرة كأساس للحركة الزمنية للخلف والأمام للاستدلال الزمنى على وقائع منذ تاريخ النشوء ووصولاً دراسات العلماء المتأنية لتطور الأنظمة البيولوجية.
- دعم التطبيقات فى مجال الطب والزراعة وغيرها من المجالات العلمية.

سيناريو A Scenario:

كمقدمة سريعة لدور طرق الحساب الآلى فى البيولوجيا الجزيئية، دعنا نتخيل أن هناك مأساة فى المستقبل حيث ظهور فيروس جديد يسبب مرضاً وبائياً للإنسان أو الحيوان. فى هذه الحالة سوف يقوم العلماء بعزل المادة الوراثية معملياً وتحديد التتابع وحينئذ تستخدم برامج الكمبيوتر.

وباستعراض هذا الجينوم الجديد ومضاهاته بالبيانات الوراثية المعروفة والموجودة فى بنك البيانات سوف يتم تعريف وتشخيص ذلك الفيروس وكذلك تحديد علاقته بالفيروسات التى سبق دراستها (١٠). تستمر بعد ذلك الدراسات بهدف انتاج علاج مضاد للفيروس. وبما أن الفيروسات تحتوى على جزيئات بروتين والتى تعتبر كأهداف مناسبة للعقاقير التى تتداخل مع تركيب ووظيفة الفيروس وأن تتابع الأحماض الأمينية للبروتينات عبارة عن رسائل مكونة من ٢٠ حرفاً من الحروف الهجائية، لذا باستخدام تتابع الحمض النووى دنا DNA يمكن بالاستعانة ببرامج الحاسب الآلى استنتاج تتابعات الحمض الأمينى فى واحد أو أكثر من بروتينات الفيروس والتى لها دوراً حاسماً فى التكرار والمضاهاة (٠١).

ومن خلال تتابع الحمض الأمينى سوف تقوم البرامج بالحساب الآلى لتراكيب تلك البروتينات على أساس أن تتابع الأحماض الأمينية يحدد التركيب ثلاثى الأبعاد لها وبالتالي خصائصها الوظيفية. فى بداية الأمر

سوف يتم مراجعة البيانات الموجودة في بنك المعلومات والخاصة بالبروتينات قريبة الصلة والمعروف تركيبها (١٥). اذا تم التعرف على أحد التراكيب فان مشكلة التنبؤ بالتركيب الجديد سوف تختزل الى أدنى درجة ممكنة (التنبؤ بالتغيرات في التتابع) وهنا يمكن التنبؤ بتركيب الأماكن المستهدفة للبروتين باستخدام نماذج التماثل Homology modeling (٢٥). أما في حالة عدم وجود تراكيب ذات علاقة ويبدو أن بروتين الفيروس جديد تماما فان عملية التنبؤ بالتركيب يجب أن تجرى بالكامل *ab initio* (٥٥). وسوف يقل في المستقبل تكرار مثل هذه الحالات حيث أنه بمرور الوقت سيحدث نمو في حجم البيانات التي يمكن تزويد البنك بها وبالتالي تزداد المقدرة على الكشف عن تراكيب عديدة.

بمعرفة تركيب بروتين الفيروس يصبح من الممكن تصميم مادة للعلاج. وحيث تمتلك البروتينات أماكن على سطح الجزيء تتناسب مع وظيفة البروتين والتي يمكن تعطيلها، فانه من الممكن تعريف وتصميم جزيء صغير متوافق من حيث الشكل والشحنة مع المكان المستهدف على سطح البروتين ليعمل كدواء مضاد للفيروس (٥٠). وهناك طريقة بديلة تعتمد على تصميم أجسام مضادة لبروتين الفيروس ثم تصنيعها واستخدامها لمعادلة الفيروس (٥٠).

يقوم هذا السيناريو على أسس ثابتة وليس هناك مجالاً للشك أنه في يوم ما سوف يطبق كما هو.

هناك سبب وحيد يحول دون تطبيق ذلك السيناريو لموجهة الفيروس المسبب لمرض الايدز وهو أن تلك الفيروسات تمتلك قدرة ما على الحماية الذاتية. وعلماء الكمبيوتر عند قراءتهم لهذا الكتاب يدركون أن الأرقام المذكورة بين أقواس في هذا الكتاب للم تستخدم للاستدلال على مراجع ولكنها تتبع نظام D. E. Knuth لفهرسة درجة صعوبة مشكلة تحت البحث كما ورد في كتابه فن برمجية الحاسب الآلي The art of computer

programming حيث تدل الأرقام أقل من ٣٠ على مشاكل ذات حلول موجودة بالفعل أما الأرقام الأعلى تدل على موضوعات قيد البحث. وأخيرا فانه يجب أن يكون من المعلوم أن الطرق التجريبية البحتة لايجاد مواد مضادة للفيروس سوف تظل ولسنوات عديدة قادمة أكثر نجاحا من التوجهات النظرية. ٦

تعتبر مادة الوراثة دنا DNA وفي بعض الفيروسات رنا RNA هي سجل المعلومات في الكائنات (طبيعة التطور والنشاط لكل فرد). وكما هو معروف فان جزيئات دنا عبارة عن سلسلة طويلة خطية تحتوي على رسالة مكونة من أربعة حروف هجائية. والرسالة طويلة حتى في الكائنات الدقيقة تتكون من ٦*١٠ حرفا تماما. ويتضمن تركيب الدنا على آليات للنسخ الذاتي وتكوين البروتينات. يؤدي التركيب الحلزوني المزدوج والمتماثل داخلها الى نسخ دقيق. ومع أن النسخ الدقيق ضروري لثبات الصفة الوراثية إلا أن بعض عمليات النسخ غير الدقيق أو آليات نقل مادة وراثية غريبة تكون أحيانا ضرورية لحدوث النشوء في الكائنات اللاجنسية.

شرائط الحلزون المزدوج تكون في وضع متوازيًا وعكسيًا ويعرف بالاتجاه ٣-٥ وذلك بالنسبة لأماكن حلقة الداى أوكسى ريبوز وعند ترجمتها الى بروتين يقرأ تتابع الحمض النووي في الاتجاه ٥-٣. ويتم تنفيذ المعلومة الوراثية من خلال تخليق الرنا والبروتينات. والبروتينات هي الجزيئات المسؤولة عن غالبية التراكيب والأنشطة في الكائنات فالشعر والعضلات والانزيمات والمستقبلات والأجسام المضادة جميعها بروتينات وكل من الأحماض النووية والبروتينات جزيئات تتكون من سلسلة طويلة خطية. تتكون الشفرة الوراثية من ثلاثة حروف متتالية من تتابع الدنا تحدد أحماض أمينية متتالية وبامتداد تتابعات الدنا يشفر تتابعات

الدوجما: مركزيا وطرفيا Dogmas: Central & Peripheral

الأحماض الأمينية في البروتينات. كذلك تتكون البروتينات من ٢٠٠ - ٤٠٠ حمض - أميني تتطلب من ٦٠٠ - ١٢٠٠ حرفا من رسائل الدنا لتحديدتها. تخليق جزيئات الرنا أيضا يتحكم فيه تتابعات الدنا بالرغم من أنه في معظم الكائنات ليس كل الدنا يترجم الى بروتينات ورنـا. بعض المناطق في تتابع الدنا يختص بالتحكم في آليات محددة كما توجد كميات كبيرة من الجينوم في كائنات راقية يبدو أنه لا فائدة منها Junk بمعنى أنه حتى الآن لم يتم التعرف أو فهم وظائفها.

في جزيئات الدنا التماثل في الحروف الهجائية يؤدي الى التشابه الكيميائي والتوحد في الشكل. وعلى العكس - تظهر البروتينات اختلافات كبيرة في التوزيع ثلاثي الأبعاد. وذلك ضروري لدعم التنوع التركيبي والوظيفي الكبيرين لها. يحدد تتابع الحمض الأميني في البروتين التركيب ثلاثي الأبعاد له. يوجد لكل تتابع طبيعي للحمض الأميني حالة فطرية ثابتة مميزة والمتأقلمة تلقائيا تحت الظروف المناسبة. عندما يتم تسخين بروتين منقى أو يتعرض لظروف مغايرة للبيئة الفسيولوجية الطبيعية ينتج تركيب غير حلزوني ذو ترتيب مختلف وغير نشط بيولوجيا. لذلك تمتلك الثدييات آليات للحفاظ على ثبات درجة الحرارة داخل أجسامها. وتستعيد جزيئات البروتين تركيبها الفطري عند العودة للظروف الطبيعية وبصورة مماثلة للجالاة الأصلية.

الطى التلقائي في جزيئات البروتينات لتكوين حالتها الفطرية هي النقطة التي عندها تقوم الطبيعة بالقفزة الضخمة من البعد الأحادي لعالم تتابعات الوراثة والبروتين الى العالم ثلاثي الأبعاد الذي نعيشه. وهنا يوجد تناقض ظاهري فترجمة تتابعات الدنا الى تتابعات حمض أميني من السهل أن توصف منطقيا بواسطة الكود الوراثي بينما من الصعب جدا الوصف المنطقي للالتفاف الدقيق للسلسلة متعددة الببتيدات الى تركيب ثلاثي الأبعاد. فبينما تتطلب عملية الترجمة الآلية الضخمة المركبة للريبوسوم

– الرنا الناقل وجزيئات مصاحبة – الا أن عملية طي البروتين تحدث تلقائيا.

تعتمد وظائف البروتينات على التركيب الثلاثي الأبعاد الفطري. على سبيل المثال يحتوى تركيب أى انزيم على تجويف على سطح الجزيء يقوم بالارتباط بجزيء صغير ليجاور أحماض أمينية حفازة. لذلك نذكر النموذج التالي:

- يحدد تتابع الدنا تتابع البروتين
- يحدد تتابع البروتين تركيب البروتين
- يحدد تركيب البروتين وظيفة البروتين

تركز معظم أنشطة المعلوماتية الحيوية المنظمة على تحليل البيانات ذات الصلة بتلك العمليات.

الى هذا الحد لم يتضمن هذا النموذج على مستويات أعلى من المستوى الجزيئى للتركيب والتنظيم وتشتمل على سبيل المثال على أسئلة كتلك المتعلقة بكيفية تخصص الأنسجة أثناء التطور أو أكثر تعميما كيف تمارس التأثيرات البيئية تحكمها فى الأحداث الوراثة. فى بعض الحالات البسيطة يكون من المفهوم على المستوى الجزيئى كيف أن زيادة كمية مادة تفاعله تسبب زيادة انتاج الأنزيم المحفز لتحوله. والأكثر تعقيدا هى برامج التطور خلال فترة حياة الكائن. تلك المشكلات الساحرة المتعلقة بتدفق المعلومات فى الكائن والتحكم فيها أصبحت الآن تأتى من خلال مجال البيومعلوماتية.

يتضمن بنك البيانات على سجل للمعلومات وتنظيم أو كيان منطقي لتلك المعلومات وأدوات للاتصال به. يغطى بنك البيانات للبيولوجيا الجزيئية تتابعات الحمض النووي والبروتين وتراكيب ووظائف الجزيئات الكبيرة. وتشتمل على:

**المشاهدات وسجلات
البيانات
Observables and
Data Archives**

- بنك بيانات سجلية للمعلومات البيولوجية:
 - تتابعات الدنا والبروتين متضمنة الشرح والتفسير.
 - تراكيب الحمض النووي والبروتين وشرحها. بنك بيانات لطرز تعبير البروتين.
 - بنك بيانات فرعية: تحتوي على المعلومات المجمعة من بنوك البيانات السجلية والناجمة عن تحليل محتواها. فعلى سبيل المثال:
 - بواعث تتابع (خصائص عائلات البروتينات)
 - الطفرات والاختلافات في تتابعات الدنا والبروتين
 - تصنيف أو علاقات (الصلات والخصائص المشتركة للمدخلات في السجلات. مثل بنك بيانات لمجموعة من عائلات تتابع البروتين أو تصنيف متسلسل لطرز النفاذ البروتين.
 - بنك بيانات مرجعية
 - بنك بيانات لمواقع الوب
 - بنك بيانات لبنوك البيانات المحتوية على معلومات بيولوجية
 - روابط بين بنوك البيانات
- تساؤلات قاعدة البيانات تبحث تعريف مجموعة من المدخلات (على سبيل المثال: تتابعات أو تراكيب) على أساس صفات محددة أو على أساس التشابه مع مجس للتتابع أو التركيب. أكثر التساؤلات شيوعا هو: عند تقدير تتابع أو تركيب جديد ما هو درجة التشابه بينه وبين الموجود في بنك البيانات؟ بمجرد التوصل الى مجموعة من التتابعات أو التراكيب من قاعدة بيانات ملائمة تتشابه مع المجس يصبح الباحث قادرا على تعريف وبحث خصائصها العامة.
- آليات الاتصال بينك للبيانات هي مجموعة من الأدوات للاجابة على الأسئلة الآتية:

- هل يحتوى بنك البيانات على المعلومات المطلوبة ؟ (مثال: فى أى من بنوك البيانات يمكن الحصول على تتابعات الحمض الأمينى للكحول ديهيدوروجينيزز؟)
- كيف يمكن جمع معلومات منتقاة من بنك البيانات فى صورة مفيدة؟ (مثال: كيف يصنف قائمة لتتابعات الجلوبيين أو جدول لمصفوفة تتابعات الجلوبيين؟)
- فهارس بنوك البيانات تفيد عند التساؤل: أين يمكن أن يوجد بعض المعلومات المحددة ؟ (مثال: ما هى بنوك البيانات التى تحتوى تتابع الحمض النووى لتريسين بوركوبيين؟) وبالطبع اذا تم معرفة وتحديد ما هو المطلوب بدقة عندئذ تبدأ خطوات حل المشكلة.

وبنك البيانات بدون طرق فعالة للاتصال والتعامل تكون بمثابة مقبرة للبيانات. كيف تحقق اتصال فعال مع قاعدة بيانات مصممة على أن تظل محجوبة عن المستفيدين. أصبح من الواضح أن الاتصال الفعال لا يمكن تزويده بنظام تساؤلات فى أرشيف غير مشبع. بدلا من ذلك يجب أن يصمم التنظيم المنطقى لتخزين المعلومات مع وجود تصور لطرق الاتصال والتعامل - ماهى نوعية الأسئلة التى يريد المستخدم طرحها - كما ينبغى أن ينسجم تركيب الأرشيف بسلسلة مع برامج استدعاء المعلومات.

تتضمن مختلف أنواع استعلام قاعدة البيانات الممكنة فى مجال البيومعلوماتية الآتى:

- (1) تتابع مفترض، أو جزء من تتابع، وإيجاد تتابعات فى قاعدة البيانات مماثلة له. هذه مشكلة مركزية فى البيومعلوماتية. يشارك المعلوماتية الحيوية مجالات عديدة من علم الحاسبات فى سلسلة من المشاكل المتماثلة. على سبيل المثال، برامج معالجة الكلمات والتحرير التى تدعم سلسلة من وظائف البحث.

(٢) تركيب بروتين مفترض، أو جزء منه ، وإيجاد تراكيب مماثلة فى قاعدة البيانات. ذلك هو التعميم لسلسلة المشاكل المتشابهة للأبعاد الثلاثة.

(٣) تتابع مفترض لبروتين غير معروف تركيبه، إيجاد تراكيب فى قاعدة البيانات والى تتبنى تراكيب ثلاثية الأبعاد مشابهة. يحث ذلك على البحث فى بنوك بيانات التتابع عن بروتينات لها تتابعات مماثلة لتتابع المجس: من المتوقع فى حالة وجود اثنان من البروتينات لهما تتابعات متماثلة لدرجة كافية فسوف يكون لهما تراكيب متشابهة. الا أن هذا القول غير حقيقى. ونأمل فى الوصول الى تقنيات بحث أكثر قوة والى ستستطيع التعرف على بروتينات متشابهة التركيب حتى ولو كان تتابع كل منهما ينحرف عن النقطة التى نقر التشابه بينهما على أساس مقارنة التتابع.

(٤) تركيب مفترض لبروتين لإيجاد تتابعات فى بنك البيانات والى تقابل تراكيب مماثلة. مرة أخرى يمكن استخدام هذا التركيب كمجس لتركيب فى بنك البيانات لكن ذلك سوف يحقق فقط نجاحا محدودا لأن هناك العديد من التتابعات المعروفة أكثر من التراكيب. ولذا فانه من المرغوب فيه وجود طريقة يمكنها التقاط التركيب من التتابع.

أرقام (١) و (٢) عبارة عن أمثلة محلولة لعمليات بحث تجرى آلاف المرات كل يوم. بينما (٣) و (٤) عبارة عن مجالات نشطة للبحث.

تتأتى أهداف غاية فى الرقة عند الرغبة فى دراسة العلاقات بين معلومات بنوك بيانات منفصلة. حيث يتطلب ذلك وجود روابط تسهل الاتصال المتزامن مع عدة بنوك بيانات. مثال على ذلك: هل يوجد فى الخميرة بروتين مماثل لذلك البروتين معلوم التركيب والذى يساهم فى حدوث أمراض تخليق البيورين فى الإنسان؟ والخلفية هنا يحددها: تركيب معلوم ووظيفة معينة واكتشاف الصلة والارتباط بالمرض ونوع معين من

الكائنات. أدى الاهتمام بالنمو المتزايد بطرق الاتصال المتزامن بينوك البيانات الى بحث التفاعل داخل بيانات البنك بمعنى - كيف يتم تبادل البيانات يبين بنك وآخر بدون تضحية كبيرة في حرية كل بنك في بناء بياناته بطرق تلائم الصفات الفردية المميزة للمادة التي تحتويها تلك البيانات.

والمشكلة التي لم تظهر بعد في مجال البيولوجيا الجزيئية هي عملية تحديث السجلات. ففي نظام قاعدة البيانات للحجز بشركات الطيران يمنع حدوث بيع المكان الواحد لأكثر من مسافر مع وجود منافذ للحجز متعددة. في مجال المعلوماتية الحيوية يمكن للمستخدم قراءة واستخلاص المعلومات من سجلات بنوك البيانات أو تقديم مواد يتم معالجتها بواسطة القائمين بالعمل في سجل ما ولكنه لا يستطيع اضافة أو تغيير المدخلات مباشرة. قد يتغير هذا الوضع. من منظور عملي تتزايد كمية البيانات الناتجة بسرعة كبيرة لدرجة تعوق قدرة مشروعات السجلات على استيعابها. ويوجد بالفعل تحرك نحو مشاركة أكبر للعلماء في المعامل لتجهيز بيانات للسجلات.

بالرغم من جدل بخصوص تحكم متميز للسجلات إلا أن هناك حاجة الى الحد من الطرق العامة (غير المتخصصة) للتعامل معها - تصميم front ends - ويمكن لمجتمعات المستخدم المتخصص أن تستخلص تحت مجموعات من البيانات أو دمج بيانات من مصادر مختلفة والوصول الى طرق متخصصة للاتصال والتعامل. وتعتمد مثل - دكاكين قواعد البيانات هذه - على السجلات الأولية كمصدر للمعلومات التي تحتويها ولكن مع اعادة تصميم التنظيم والعرض بالطريقة التي يرونها أكثر ملاءمة. وفي الحقيقة يمكن لمختلف قواعد البيانات الفرعية أن تصنف ذات المعلومة. يقترح الاستقراء المعقول مفهوم قواعد البيانات الواقعية المتخصصة virtual data bases المبنى على السجلات وفي نفس الوقت يقدم مفهوما ووظيفة فردية فصلت لتلبي احتياجات مجموعات البحث الفردية أو حتى علماء منفردين.

الرعاية والتفسير والتحكم في الجودة:

Curation Annotation, and Quality Control:

تعتمد المجتمعات العلمية والطبية على جودة بنوك البيانات. ومؤشرات الجودة حتى وان كانت لا تسمح بتصحيح الأخطاء الا أنها قد تتجنب التوصل الى استنتاجات خاطئة.

تضمن مدخلات بنك البيانات النتائج التجريبية الخام ومعلومات معاونة أو تفسيرات. وتمتلك كل واحدة من تلك مصادرها الخاصة من الخطأ.

من أهم المحددات لجودة البيانات ذاتها هو مدى دقة التجارب. البيانات القديمة محدودة بتقنيات قديمة. فعلى سبيل المثال: تتابعات الحمض الأميني للبروتينات التي تم تقديرها بواسطة تحديد التتابع البيتيدي يتم الآن ترجمتها جميعا من تتابعات الدنا. أحد عواقب انفجار البيانات هو أن غالبية البيانات جديدة وتم الحصول عليها بتكنولوجيا حديثة والتي في معظم الحالات تؤدي عملا جيدا.

تشتمل التفسيرات معلومات حول مصادر البيانات والطرق المستخدمة في تقديرها. كما أنها تقدم روابط مع المعلومات ذات الصلة في بنوك المعلومات الأخرى. في بنوك بيانات التتابع تتضمن التفسيرات جداول توصيف: قوائم بأجزاء التتابعات التي لها معنوية بيولوجية - على سبيل المثال مناطق من تتابع الدنا الخاصة بشفرة البروتينات. يظهر ذلك في تصميم مناسب للتعامل مع الحاسب الآلي ويكون محتوياتها مقيدة بمفردات لغوية محكمة. حتى وقت قريب ادخال تتابع نمطي لدنا يتم بواسطة مجموعة بحثية منفردة تبحث الجين ونواتجه بطريقة مترابطة منطقيا. تستند التفسيرات على البيانات التجريبية والمكتوبة بواسطة متخصصون. على النقيض، لا تقدم مشروعات تتابع الجينوم تأكيد تجريبي للتعبير في معظم الجينات المفترضة ولا تشخيص نواتجها. يقيم الأمناء في بنوك البيانات تفسيراتهم على أساس تحليل التتابعات بواسطة برامج الحاسب الآلي.

التفسيرات هي أضعف مكون في مؤسسة الجينوم ويمكنه التفسيرات تكون ممكنة فقط بدرجة محدودة. وللحصول عليها بطريقة صحيحة تكون كثيفة

العمالة وتقسيم المصادر يكون غير كافيا. ولكن لا يمكن الاستخفاف بالتفسيرات الصحيحة. وقد علق بورك بأن تلك الأخطاء في دراسة الجين تفسد جودة بيانات التتابع. سوف يؤدي نمو بيانات الجينوم الى تحسين جودة التفسير كما أن الطرق الاحصائية تزيد من الدقة وتسمح باعادة محسنة لتفسير المدخلات. التحسن في التفسيرات شيء جيد ولكن التلازم الحتمي - التفسير سيكون متدفقا - سيكون مزعجا. هل يجب الاطلاع دوريا على البحوث المكتملة والأخذ في الاعتبار خلاصة ما توصلت له من نتائج؟ وقد تفاقمت المشكلة بوجود العديد من مواقع الوب مع تزايد شبكات الربط الكثيفة. ويتيح ذلك طرق مفيدة للتطبيقات المتنوعة. والوب أيضا ناقل لعدوى تضخم الأخطاء في البيانات الخام والتفسيرات المختلفة ويتم تباعا تصويب الأخطاء في البيانات في صورتها الأولية ولكن الحاجة الى التصويب لا نهاية لها.

والحل الوحيد الممكن هو تصويب موزع وديناميكي لمعالجة الخطأ والتفسير. سوف يقوم الأخصائيون الموزعون بدور الأمناء حيث أن العاملين بينك البيانات ليس لديهم لا الوقت ولا الخبرة للقيام بتلك المهمة. كما تسمح ديناميكية التقدم في ميكنة تعريف وتصحيح الخطأ والتفسير باعادة التفسير لبنوك البيانات. وسوف نضطر أن نسلم بالفكرة الأمانة لبنك بيانات مستقر يتكون من مدخلات سليمة من بداية توزعها وتظل كذلك. وسوف تصبح بنوك البيانات مصدر غال وثمين للمعلومات والتي تنمو في الحجم وتزداد نضجا ونأمل أن تكون بجودة مناسبة.

تستخدم الشبكة العنكبوتية للحصول على مواد مرجعية و أخبار وللتعامل مع قواعد البيانات في البيولوجيا الجزيئية أو لمجرد التصفح. وتعتبر الوب الآن وسيلة من وسائل الاتصال بين الأشخاص وبين الحاسبات عبر الشبكات. وهي بمثابة قرية عالمية تشتمل على ما يقابل المكتبة ومكتب البريد والأسواق والمدارس وغيرها.

**الشبكة العنكبوتية
العالمية واسعة الانتشار
The World Wide
:Web (www)**

ويجرى المستخدم برنامج للتصفح على الحاسب الخاص به. ومن برامج التصفح الشائعة:

Netscape and Internet Explorer ويمكن بواسطة تلك البرامج قراءة وعرض مواد من أي مكان في العالم. تقدم برامج التصفح معلومات للتحكم كما تتيح نقل المعلومات إلى حاسب محلي. ومن الأشياء الرئيسية لبدء استخدام الشبكة بفاعلية إيجاد نقاط دخول مفيدة حيث تأخذك روابط الوصل (links) إلى المكان الذي تريده بمجرد بدء التشغيل. ومن بين أكثر المواقع أهمية تلك الخاصة ببرامج البحث والتنقيب search engines والتي تفهرس الوب بأكمله وتتيح استرجاع المعلومات باستخدام كلمات رئيسية key words. ومن الممكن إدخال واحدا أو أكثر من المصطلحات مثل 'phosphorylase', 'allosteric change', 'crystal structure' وسوف يظهر برنامج البحث قائمة بروابط الوصل لمواقع على الوب تحتوي على تلك المصطلحات وعندئذ يمكن التعرف على المواقع محل الاهتمام.

أثناء جلسة تصفح الوب يمكن الاحتفاظ بالمستندات التي نحتاج الرجوع إليها في جلسات تصفح قادمة وذلك بحفظ روابط الوصل الخاصة بها في ملف bookmarks أو favourites ومن ثم يمكنك في جلسات لاحقة الرجوع إلى أي موقع مباشرة دون الحاجة إلى اتباع محاولات روابط الوصل المستخدمة في أول مرة.

الوب ليس طريقا ذو اتجاه واحد بمعنى أن العديد من مستندات الوب تتضمن نماذج يمكن ادخال معلومات فيها وإجراء برنامج للحصول على النتائج خلال نفس الجلسة. وتعتبر برامج البحث والتنقيب خير مثال على ذلك. وتتطلب الآن العديد من العمليات الحسابية في البيومعلوماتية عبر أجهزة الخدمة servers. وفي حالة العمليات الحسابية الطويلة ربما لا تتأتى النتائج أثناء نفس الجلسة ولكن ترسل بالبريد الإلكتروني e-mail.

محددات مواقع المصدر The hURLy-bURLy (Uniform Resource :Locators)

تحدد تلك الحروف شكل المادة ومكانها حيث ينبغي أن يكون لكل مستند على الوب ملف في مكان ما على حاسب معين. مثال لـ URL:
<http://www.lib.berkeley.edu/TeachingLib/Guides/Internet/FindInfo.html>
 وهو موقع لدرس حول كيفية الحصول على معلومات من الإنترنت. ويشير مختصر http:// الى أن المستند في صيغة بروتوكول نقل النص المحوري hypertext transfer protocol أما www.lib.berkeley.edu فهو اسم الحاسب: المكتبة المركزية في جامعة كاليفورنيا بيركلي. وباقي الأجزاء تشير الى مكان واسم الملف في الحاسب.

النشر الإلكتروني Electronic publication

هناك العديد والعديد من المواد المنشورة على صفحات الوب. وفي المجالات العلمية قد ينشر جدول المحتويات فقط أو جدول المحتويات مع ملخصات المقالات أو المقالات كاملة. والآن يظهر العديد من المطبوعات المؤسسية والنشرات الدورية والتقارير الفنية على الوب. كما يحتوى الكثير من المطبوعات على مراجع لروابط وصل تحتوى على مواد مساعدة لا يمكن ظهورها على ورق. كما أن المواد المطبوعة على ورق يمكن أن تتضمن عناوين لمواقع على الوب وللبريد الإلكتروني.

الحاسبات وعلم الحاسب Computers and computer science

لم يكن من المحتمل ظهور مجال البيومعلوماتية بدون التقدم الذى تحقق فى المكونات المادية Hardware وبرامج software الحاسبات. كما أن وسائط الحفظ السريعة ذات الكفاءة العالية ضرورية للإبقاء على السجلات. وتتطلب عملية استرجاع وتحليل المعلومات برامج بعضها بسيط ولس والبعض الآخر متطورة ومعقدة. كما يتطلب توزيع

المعلومات توافر إمكانيات من شبكات الحاسب computer networks وكذلك بالنسبة للشبكة العنكبوتية العالمية. www.

علوم الحاسب مجال صغير ومزدهر يستهدف تعظيم الاستخدام الفعال لمكونات تكنولوجيا المعلومات. مناطق معينة من علم الحاسب تمس البيومعلوماتية مسا وثيقا ومباشرا. دعنا نتصور مشكلة بيولوجية مثل استرجاع كل التتابعات المماثلة لمجس التتابع probe sequence من قاعدة بيانات سيكون الحل الصواب هو اللجوء الى علم الحاسب في:

• التحليل العددي (الخوارزميات) Analysis of algorithms:

وهو الوصف الكامل والدقيق لطريقة حل المشكلة. بالنسبة لاسترجاع تتابعات مماثلة نحتاج قياس مدى تشابه تتابع المجس لكل تتابع موجود في قاعدة البيانات. ومن المحتمل أن يكون ذلك أفضل بكثير من تلك الاتجاه الساذج لفحص كل زوج من الأماكن في كل تجاور محتمل وهي طريقة حتى بدون السماح بفراغات فانها تحتاج الى وقت يتناسب مع ناتج ضرب عدد حروف مجس التتابع في عدد حروف تتابعات قاعدة البيانات. ويركز تخصص الحاسب والمعروف عامة 'stringology' على ايجاد طرقا ذات كفاءة للتعامل مع هذا النوع من المشاكل وتحليل فاعلية أدائها.

• تراكيب البيانات واسترجاع المعلومات Data structures, and information retrieval:

كيف يتم تنظيم البيانات بطريقة تتيح استجابة كفاء للتساؤلات؟. على سبيل المثال: هل هناك طرق لفهرسة أو اعداد معالجة للبيانات لجعل بحث تماثل التتابع أكثر كفاءة؟ كيف يمكننا تقديم حدودا مشتركة من شأنها مساعدة المستخدم على تصميم و تنفيذ التساؤلات؟.

• هندسة البرامج Software engineering:

لم يعد من الصعب بتاتا كتابة برامج بلغة الحاسبات الأصلية. يعمل المتخصصون في عمل البرامج بلغات عالية المستوى مثل

C, C++, PERL ('Practical Extraction and Report Language')

أو حتى لغة FORTRAN. يعتمد اختيار لغة البرمجة على طبيعة الحساب وتركيب البيانات المصاحبة. وبالطبع فمعظم البرامج المعقدة المستخدمة في البيومعلوماتية تكتب بواسطة متخصصون.

البرمجة Programming:

البرمجة بالنسبة لعلم الحاسب كالقائم بالبناء في فن العمارة كلاهما مبدع أحدهما فن والآخر حرفة.

يستفسر العديد من طلاب البيومعلوماتية هل من الضروري أن يتعلموا كتابة برامج حاسب معقدة؟ والاجابة أنه ليس من الضروري ذلك إلا إذا كانت هناك رغبة في التخصص في هذا المجال. ويتطلب العمل في مجال البيومعلوماتية اكتساب خبرات في استخدام الأدوات المتاحة على صفحات الوب. كما أنه من الأشياء الأساسية هو تعلم كيفية إنشاء موقعا على الوب وكذلك الإبقاء عليه وبالطبع هناك حاجة الى إمكانيات لاستخدام نظام تشغيل الحاسب الشخصي. ومهارة كتابة نصوص بسيطة بلغة مثل PERL تعتبر من ضمن أساسيات نظام التشغيل.

وحيث يجب أن يؤخذ في الاعتبار حجم سجلات البيانات والنمو المتزايد في درجة التعقيد في التساؤلات المطروحة لذلك من الأفضل أن يترك الابداع الحقيقي للبرمجة في هذا المجال للمتخصصين ذى الخبرة الجيدة في علم الحاسب.

وينصح بتعلم المهارات الأساسية للغة PERL لأنها أداة قوية تجعل من السهل جدا القيام بالعديد من العمليات البسيطة والمفيدة. وتمتاز لغة PERL بأنها متاحة في غالبية أنظمة الحاسب.

كيف يمكنك تعلم PERL بدرجة كافية لاستخدامها في البيومعلوماتية؟ تقدم العديد من المعاهد دروسا لهذه اللغة كما يمكن تعلمها بمساعدة

الزملاء وكذلك بالرجوع الى الكتب المتوفرة. ومن الطرق المفيدة أيضا البحث عن دروس على صفحات الوب بالاستعانة ببرامج البحث حيث يوجد مواقع لذلك. يمكن الرجوع الى موقع مشروع بيوبيرل Bioperl project : (<http://bio.perl.org/>) والذي يتيح مصدرا لبرامج PERL ومكوناته المستخدمة في مجال البيومعلوماتية.

قوة PERL في تناول سلسلة الحروف جعلها تلائم عمليات تحليل التتابع في علم البيولوجي. وفيما يلي مثال لاستخدام برنامج PERL بسيط لترجمة تتابع نيكليوتيدة الى تتابع حمض نووي طبقا لكود وراثي قياسي. السطر الأول: `#!/usr/bin/perl` (هو اشارة الى نظام تشغيل UNIX (or LINUX) وما يليه هو برنامج PERL. خلال البرنامج جميع النصوص التي تبدأ ب `#` ليست الا تعليق. ويشير السطر `_END_` الى انتهاء البرنامج وأن ما يليه هو عبارة عن بيانات تم ادخالها. (يمكن الحصول على المواد والبرامج بالرجوع الى موقع الكتاب على الوب وهو:

<http://www.oup.com/uk/lesk/bioinf>

ويعرض هذا المثال البسيط صور عديدة للغة PERL. ويحتوى هذا الملف على بيانات مصاحبة (جدول ترجمة الكود الوراثة)، وعبارات تخبر الحاسب لعمل شئ معين بالمدخلات (مثل التتابع المطلوب ترجمته)، والبيانات التي يتم ادخالها (وهي تظهر بعد سطر `_END_`). كما تلخص التعليقات أقسام من البرنامج وتوصف تأثير كل عبارة.

يتركب البرنامج كبلوكات داخل أقواس متعرجة: { ... } مما يفيد في انسياب الأداء. وداخل البلوكات عبارات فردية (كل منها تنتهى ب ;). والبلوك الخارجى عبارة عن لوب:

```
while ($line = <DATA>) {
```

```
...
```

```
}
```

تشير <DATA> الى سطور ادخال البيانات (والتي تظهر بعد _END_). ويتم اجراء البلوك مرة واحدة لكل سطر من المدخلات ويستمر ذلك حتى نهاية كل السطور.

ويظهر في البرنامج ثلاثة أنواع من تراكيب البيانات. سطر ادخال البيانات ويشار اليه \$ line وهو سلسلة من الحروف البسيطة والتي تجزأ الى منظومة متعددة البيانات array أو حامل لثلاثيات triplets. وتخزن المنظومة بيانات لموضوعات عديدة في ترتيب خطي. ويمكن استرجاع بيانات كل موضوع على حدي من أماكنها في المنظومة. لتسهيل التقاط الكود الثلاثي لحامض أميني فإنه يتم تخزين الكود الوراثةي كمنظومة ترتيب. ومنظومة الترتيب أو جدول التكرار عبارة عن تعميم لمنظومة بسيطة أو تسلسلية. اذا كانت عناصر المنظومة البسيطة مفهومة بواسطة أرقام متتابعة فان عناصر منظومة الترتيب تفهرس باستخدام سلاسل من حروف وهي في هذه الحالة تكون 64 ثلاثية. يتم معالجة الثلاثيات المدخلة طبقاً لترتيب ظهورهم في تتابع النيكلويدية مع الاستعانة بعناصر جدول الكود الوراثةي بترتيب تحملي كما يقرأ في الثلاثيات المتتالية. المنظومة البسيطة أو حامل سلاسل الحروف يكون مناسباً لمعالجة ثلاثيات متتالية بينما تلائم المنظومة التسلسلية التقاط الأحماض الأمينية المقابلة للثلاثيات.

مثال : برنامج بيرل لترجمة تتابع حامض نووي الى تتابع حامض أميني

Translate.pl - PERL program to translate nucleic acid sequence to amino acid sequence:

```
#!/usr/bin/perl
#translate.pl -- translate nucleic acid sequence to protein
                sequence
#
                according to standard genetic code
#
# set up table of standard genetic code

%standardgeneticcode = (
  "ttt"=> "Phe",   "tct"=> "Ser", "tat"=> "Tyr",   "tgt"=> "Cys",
  "ttc"=> "Phe",   "tcc"=> "Ser", "tac"=> "Tyr",   "tgc"=> "Cys",
  "tta"=> "Leu",   "tca"=> "Ser", "taa"=> "TER",   "tga"=> "TER",
  "ttg"=> "Leu",   "tcg"=> "Ser", "tag"=> "TER",   "tgg"=> "Trp",
  "ctt"=> "Leu",   "cct"=> "Pro", "cat"=> "His",   "cgt"=> "Arg",
  "ctc"=> "Leu",   "ccc"=> "Pro", "cac"=> "His",   "cgc"=> "Arg",
```

```

"cta"=> "Leu", "cca"=> "Pro", "caa"=> "Gln", "cga"=> "Arg",
"ctg"=> "Leu", "ccg"=> "Pro", "cag"=> "Gln", "cgg"=> "Arg",
'att'=> "Ile", 'act'=> "Thr", 'aat'=> "Asn", 'agt'=> "Ser",
'atc'=> "Ile", 'acc'=> "Thr", 'aac'=> "Asn", 'agc'=> "Ser",
'ata'=> "Ile", 'aca'=> "Thr", 'aaa'=> "Lys", 'aga'=> "Arg",
'atg'=> "Met", 'acg'=> "Thr", 'aag'=> "Lys", 'agg'=> "Arg",
'gtt'=> "Val", 'gct'=> "Ala", 'gat'=> "Asp", 'ggt'=> "Gly",
'gtc'=> "Val", 'gcc'=> "Ala", 'gac'=> "Asp", 'ggc'=> "Gly",
'gta'=> "Val", 'gca'=> "Ala", 'gaa'=> "Glu", 'gga'=> "Gly",
'gtg'=> "Val", 'gcg'=> "Ala", 'gag'=> "Glu", 'ggg'=> "Gly"
}

# process input data

while ($line = <DATA>) { # read in
line of input #
print "$line"; #
transcribe to output
chop(); # remove
end-of-line character #
@triplets = unpack("a3" x (length($line)/3), $line); # pull out
successive triplets #
foreach $codon (@triplets) { # loop
over triplets #
print "$standardgeneticcode($codon)"; # print
out translation of each #
} # end loop
cn triplets # skip
print "\n\n"; #
line on output # end loop
}
cn input lines

# what follows is input data

__END__
atgcatccctttaa
tctgtctga

```

Assemble.pl - PERL program to assemble overlapping fragments of strings:

```

#!/usr/bin/perl
#assemble.pl -- assemble overlapping fragments of strings

# input of fragments
while ($line = <DATA>) { # read in fragments, 1
per line #
chop($line); # remove trailing
carriage return #
push(@fragments,$line); # copy each fragment
into array #
}
# now array @fragments contains fragments

# we need two relationships between fragments:
# (1) which fragment shares no prefix with suffix of another
fragment
# * This tells us which fragment comes first

```



```

# (2) which fragment shares longest suffix with a prefix of
another
# * This tells us which fragment follows any fragment

# First set array of prefixes to the default value
"noprefixfound".
# Later, change this default value when a prefix is found.
# The one fragment that retains the default value must be come
first.

# Then loop over pairs of fragments to determine maximal overlap.
# This determines successor of each fragment
# Note in passing that if a fragment has a successor then the
# successor must have a prefix

foreach $i (@fragments) { # initially set prefix
of each fragment # to
  $prefix{$i} = "noprefixfound"; #
"noprefixfound" # this will be
} # overwritten when a prefix is found

# for each pair, find longest overlap of suffix of one with prefix
of the other
# This tells us which fragment FOLLOWS any fragment

foreach $i (@fragments) { # loop over fragments
  $longestsuffix = ""; # initialize longest
suffix to null

  foreach $j (@fragments) { # loop over fragment
pairs
  unless ($i eq $j) { # don't check fragment
against itself

    $combine = $i . "XXX" . $j; # concatenate fragments,
with fence XXX
    $combine =~ /([\S ]{2,})XXX\1/; # check for
repeated sequence
    if (length($1) > length($longestsuffix)) { # keep
longest overlap
      $longestsuffix = $1; # retain longest suffix
      $successor{$i} = $j; # record that $j follows
    }
  }
}

$prefix{$successor{$i}} = "found"; # if $j follows $i then
$j must have a prefix
}

foreach (@fragments) { # find fragment that has
no prefix; that's the start
  if ($prefix{$_} eq "noprefixfound") {$outstring = $_;}
}

```

```

$test = $outstring;          # start with fragment
without prefix
while ($successor($test)) {  # append fragments in
order
    $test = $successor($test); # choose next fragment
    $outstring = $outstring."XXX". $test; # append to string
    $outstring =~ s/{[\S ]+}XXX\1/\1/; # remove overlapping
segment
}

$outstring =~ s/\\n\\n/g;    # change signal \n to
real carriage return
print "$outstring\n";      # print final result

```

__END__

```

the men and women merely players;\n
one man in his time
All the world's
their entrances,\nand one man
stage,\nAnd all the men and women
They have their exits and their entrances,\n
world's a stage,\nAnd all
their entrances,\nand one man
in his time plays many parts.
merely players;\nThey have

```

وهناك اصدار بديل من البرنامج لمضاهاة الأجزاء المتداخلة assemble : overlapping fragments

```

# /usr/bin/perl

$. = "";
@fragments = split("\n",<DATA>);

foreach (@fragments) { $firstfragment($_) = $_; }

foreach $i (@fragments) {
    foreach $j (@fragments) { unless ($i eq $j) {
        ($combine = $i . "XXX" . $j) =~ s/{[\S ]{2,}}XXX\1/;
        (length($1) <= length($successor($i))) || { $successor($i)
= $j };
    }
    undef $firstfragment($successor($i));
}

$outstring = $outstring = join "", values(%firstfragment);
while ($test = $successor($test)) { ($outstring .= "XXX" . $test)
= s/{[\S ]+}XXX\1/\1/; }

$outstring =~ s/\\n\\n/g; print "$outstring\n";

__END__
the men and women merely players;\n
one man in his time

```

All the world's
their entrances, \nand one man
stage, \nAnd all the men and women
They have their exits and their entrances, \n
world's a stage, \nAnd all
their entrances, \nand one man
in his time plays many parts.
merely players; \nThey have

تعتمد التسمية البيولوجية على تقسيم لكائنات الحية الى مملكة وقبيلة
وأقسام وأجناس وأنواع وذلك على أساس أوجه التشابه المشاهدة. وتقدم
نتائج تحليل التتابع دلالات قاطعة على العلاقات بين الأنواع.

التصنيف والتسمية البيولوجية

Biological Classification and nomenclature:

استخدام التتابع لتحديد علاقات القرابة: Use of sequence to determine
phylogenetic

توضح الأمثلة التالية تطبيقات استرجاع التتابعات من بنوك المعلومات
ومقارنات التتابع في تحليل العلاقات البيولوجية:

المثال الأول: استرجاع تتابع الحمض الأميني لانزيم الريبونوكليز في
بنكرياس الحصان باستخدام ExPASy server في المعهد السويسري
للبيومعلوماتية وعنوانه:

<http://www.expasy.ch/cgi-bin/sport-search-ful>.

١- اكتب في المكان المخصص ل key words : horse pancreatic
ribonuclease

٢- ثم اضغط على مفتاح ENTER

٣- اختار RNP_HORSE ثم FASTA format سوف تحصل على التالي:

```
>sp|P00674|RNP_HORSE RIBONUCLEASE PANCREATIC (EC 3.1.27.5) (RNASE
1) ...
KESPAMKFERQHMDSGSTSSSNPTYCNQMMKRRNMTQGWCCKPVNTFVHEP
LADVQAICLQKNITCKNGQSNCYQSSSMHITDCRLTSGSKYPNCAYQTS
QKERHIIIVACEGNPYVPVHFDASVEVST
```

وهنا يمكنك أن تقوم بعملية قص ولصق الى برامج أخرى حيث يمكن
استرجاع عدة تتابعات وعمل sequence alignment وهذا يفيد في تقدير
درجة القرابة والعلاقات.

المثال الثاني: حدد باستخدام تتابعات انزيم الريبونيوكلسيز البكترياسي للحصان والحوث والكنجر النوعين الأكثر قرابة:

Pancreatic ribonuclease sequences from horse (*Equus caballus*), minke whale (*Balaenoptera acutorostrata*) and red kangaroo (*Macropus rufus*):

```
>RNP_HORSE
KESFAMKFERQHMDSGSTSSSNPTYCNQMMKRRNMTQGWCKPVNTFVHEP
LADVQAICLQKNITCKNGQSNCYQSSSMHITDCRLTSGSKYPNCAYQTS
QKERHIIIVACEGNPYVPVHFDASVEVST
>RNP_BALAC
RESFAMKFQRQHMDSGNSPGNNPNYCNQMMRRKMTQGRCKPVNTFVHES
LEDVKAVCSQKNVLCCKNGRTNICYBSNSTMHITDCRQTGSSKYPNCAKYTS
QKEKHHIIVACEGNPYVPVHFDNSV
>RNP_MACRU
ETPAEKFQRQHMDTEHSTASSNYCNLMMKARDMTSGRCKPLNTFIHEPK
SVVDAVCHQENVTCCKNGRTNICYKSNRSLITNCRQTGASKYPNCQYETSN
LNKQIIIVACEGQYVPVHFDAYV
```

يستخدم برنامج multiple - sequence alignment program CLUSTAL W وعنوانه:

<http://www.ebi.ac.uk/clustalw/>

وهناك بديل آخر وهو T-coffee:

<http://www.ch.embnet.org/software/TCoffee.html>

وسوف تصل الى تطابق مواقع في تتابعات كلا من الحصان والحوث.

من أشهر الأمثلة هو البحث عن قاعدة بيانات لموضوعات تتشابه مع مجس. ففي حالة تحديد تتابع جين جديد أو التعرف من خلال الجينوم البشري على جين مسئول عن مرض معين يكون هناك رغبة لمعرفة وجود جينات مماثلة في أنواع أخرى. والطريقة المثلى لتحقيق ذلك يجب أن تكون حساسة بحيث تتعرف على كل العلاقات واختيارية بأن تكون تلك العلاقات حقيقية.

وتشتمل طرق البحث في قواعد البيانات المتساوب بين الحساسية والاختيارية. هل تستطيع الطريقة إيجاد كل أو معظم التطابقات الموجودة بالفعل أم أنها تفقد أجزاء كبيرة؟. وعلى النقيض كم من التطابقات الواردة

**البحث عن التتابعات
المتماثلة باستخدام قواعد
بيانات PSI-BLAST:**

تكون غير صحيحة؟. بافتراض أن قاعدة بيانات تحتوى على ١٠٠ تتابع جلوبين وأن عملية البحث فى تلك القاعدة للجلوبين أعطت ٩٠٠ تتابع، وكان منها ٧٠٠ جلوبين حقيقى و ٢٠٠ خطأ. اذا يمكن القول أن هذه النتائج تحتوى على ٣٠٠ مفقودة (نتيجة خادعة سلبية) و ٢٠٠ نتيجة خادعة ايجابية. وسينتج عن تخفيض الحد الحرج للأمان زيادة لكلا النوعين. وهنا يكون الحرص على العمل بحدود حرجة منخفضة للتأكد من عدم فقد أى شئ، ويتطلب ذلك فحص دقيق للنتائج للتخلص من النتائج الخادعة الإيجابية.

من الأدوات الفاعلة للبحث فى قواعد بيانات التتابع بالاستعانة بمجس للتتابع استخدام PSI-BLAST

(Position Sensitive Iterated - Basic Linear Alignment Sequence Tool)

وهو برنامج من المركز القومى الأمريكى لمعلومات التكنولوجيا الحيوية (NCBI). ويعمل برنامج BLAST على التعرف على مناطق التشابه بدون فراغات ثم ضمها معا. ويشير PSI الى تحسين وتهذيب نمط التعرف داخل التتابع فى المراحل الأولية للبحث فى قاعدة البيانات. يؤدى اقرار انماط متحفظة الى تعظيم كل من حساسية واختيارية البحث. ويتضمن PSI-BLAST عمليات متكررة حيث يتحسن تعريف الأنماط الناشئة من خلال المراحل المتعاقبة للبحث.

مثال: تماثل الجين البشرى PAX-6 وهى جينات تتحكم فى تطور العين فى أنواع عديدة من الكائنات. وقد وجد تماثل هذا الجين فى الانسان وذبابة الدروسفيلا.

ويمكن اجراء البحث عن التماثل كما يلى:

١- الحصول على تتابع الحمض الأمينى للبروتين بالرجوع الى SWISS-PROT entry P26367.

٢- اجراء PSI-BLAST من خلال الموقع:

<http://www.ncbi.nlm.nih.gov/blast/psiblast.cgi>.

٣- ادخل التتابع واستخدم الخيارات لتحديد قاعدة بيانات للبحث وجدول التشابه المستخدم.

سوف يظهر البرنامج قائمة مشابهة لتتابع المجس ومرتبّة تنازلياً طبقاً لدرجة المعنوية الاحصائية. وفيما يلي نموذج طبق الأصل لأحد الأسطر في القائمة:

Pir 11 I 45557 eyeless, long form - fruit fly (Drosophila melano
255 7e-67

حيث Pir عبارة عن مصدر تعريف البروتين وهو 145557 entry وهو التماثل لعين الدروسفيليا. والرقم ٢٥٥ مقياس لدرجة التطابق. و 7e-67 تدل على مدى معنوية التطابق.

كما يمكن استخلاص أسماء الأنواع من نتائج PSI-BLAST وذلك باستخدام برنامج PERL كما يلي:

PERL program for extraction of species names from PSI-BLAST output:

```
#!/usr/bin/perl
#extract species from psiblast output

# Method:
#   For each line of input, check for a pattern of form [Drosophila
#   melanogaster]
#   Use each pattern found as the index in an associative array
#   The value corresponding to this index is irrelevant
#   By using an associative array, subsequent instances of the same
#   species will overwrite the first instance, keeping only a
#   unique set
#   After processing of input complete, sort results and print.

while (<>) {
    # read line of input
    if (/^\[[A-Z]{a-z}+[a-z]+\]\//) { # select lines containing
        strings of form
        #
        # [Drosophila
        # melanogaster]
        $species{$1} = 1; # make or overwrite entry
    }
    # associative
    array
}

foreach (sort(keys(%species))) {
    # in alphabetical order,
    print "$_\n"; # print species names
}
}
```

وقد وجد أن هناك تماثل مع ٥٢ نوعاً.

التنبؤ بتركيب البروتينات وهندستها: Protein structure and engineering

يحدد تتابع الحمض الأميني لبروتين ما التركيب ثلاثي الأبعاد له. باحتواء تتابعات الحمض الأميني على معلومات كافية لتحديد التركيب ثلاثية الأبعاد للبروتينات يصبح من الممكن استنباط نظام حسابي للتنبؤ بتركيب تركيب بروتين من تتابع الحمض الأميني. بالإضافة الى التنبؤ بالتركيب فقد حدد العلماء عددا من الأهداف أقل طموحا يمكن تحقيقها:

١- التنبؤ بالتركيب الثانوي Secondary structure prediction:

تحديد ما هي أجزاء التتابع التي تكون الشكل اللولبي والأخرى التي تكون الشرائط.

٢- تمييز الطي Fold recognition:

بعمل مكتبة لتركيب بروتينات معروفة ولتتابعات الحمض الأميني لها هل يصبح في استطاعتنا الحصول من تلك المكتبة على تركيب يشابه لدرجة كبيرة نظام طي للبروتين المطلوب معرفة تركيبه؟

٣- نماذج التماثل Homology modeling:

نفترض أن هناك بروتين معروف تتابع الحمض الأميني له وغير معروف التركيب ولكنه متماثل مع واحد أو أكثر من البروتينات المعروفة تركيبها. هنا يمكننا أن نتوقع أن البروتين الأكثر تماثلا يمكن استخدامه كأساس لنموذج للبروتين المجهول التركيب. ويعتمد كمال ودقة النتائج على مدى تشابه التتابع. وعموما وجد انه في حالة ما اذا كان في اثنين من البروتينات قريبة الصلة تطابق ٥٠% في التتابع عند عمل المحازاة يكون هناك تماثل ٩٠% في تركيبهما.

فيما يلي التتابعات المصفوفة والتركيب المتطابقة لاثنين من البروتينات قريبة الصلة وهما ليسوزيم الأبيض في بيض الدجاج وألفا لاكتوالبومين

في البابون. والتتابعات فيهما شديدة التقارب (٣٧% تطابق كامل في التتابعات المتراسة) كما أن هناك تماثل في التركيب. وكل بروتين يمكن استخدامه كنموذج جيد للبروتين الآخر:

Chicken lysozyme :

KVFGRCELAAAMKRHGLDNYRGYSLGNWVCAAKFESNFNTQATNRNT
DGS

Baboon ? - lactalbumin :

KQFTKCELSQNLV- -DIDGYGRIALPELICTMFHTSGYDTQAIVEND - ES

Chicken lysozyme :

TDYGILQINSRWWCNDGRTPGSRNLCNIPCSALLSSDITASVNCAKKIVS

Baboon ? - lactalbumin :

TEYGLFQISNALWCKSSQSPQSRNICDITCDKFLDDDDITDDIMCAKKILD

Chicken lysozyme :

DGN- GMNAWVAVRNRCRGTDVQA- WIRGCR-

Baboon ? - lactalbumin :

I - - KGIDYWIAHKALC - TEKL - EQWL - - CE - K



التقويم الحرج للتنبؤ بالتركيب :Structure Prediction (CASP)

يتطلب الحكم على تقانات التنبؤ بتركيب البروتينات الى استخدام اختبارات مصمته blind tests. لهذا الغرض صمم جاي مولت برامج CASP. العلماء المنشغلون بتقدير تركيب البروتين باستخدام طرق القياس البللورى والرنين النووى المغناطيسى مدعوون الى: (١) نشر نتاج الحمض الأمينى عدة شهور قبل التاريخ المتوقع الى اكتمال التجارب، (٢) الاحتفاظ بسرية النتائج حتى تاريخ متفق عليه. يقدم القائمون بالتنبؤ بنماذج تنبؤ يحتفظ بها حتى التاريخ المتفق عليه لاعلان النتائج التجريبية. وهنا يتم مقارنة نتائج التنبؤ والتجريب. ولقد سجلت نتائج التقييم باستخدام برامج CASP تقدما فى فاعلية التنبؤ والذى يرجع الى نمو بنوك المعلومات وأيضا بسبب التحسن فى الطرق المستخدمة. والباب الخامس من الكتاب يناقش التنبؤ بتركيب البروتين.

هندسة البروتينات :Protein engineering

كان هناك تشابه بين طبيعة عمل كلا من علماء البيولوجيا الجزيئية وعلماء الفلك من حيث القدرة على ملاحظة الأشياء دون اجراء تعديلات عليها. الا أن هذا لم يعد حقيقيا الآن. ففي المعمل يمكننا عمل تعديل فى الأحماض النووية والبروتين واحدات طفرات لمعرفة التأثير على وظائفهم. فمن الممكن أن يصبح لبروتين قديم وظيفة جديدة مثل انتاج أجسام مضادة حفازة ومحاولة تحضير بروتينات جديدة.

تشتق العديد من قواعد تركيب البروتينات من ملاحظات للبروتين الطبيعى. وليس ضروريا تطبيق تلك القواعد فى البروتينات المهندسة. تمتلك البروتينات الطبيعية صفات مرتبطة بالأسس العامة للكيمياء الطبيعية وآليات نشوء البروتين. ويجب أن تخضع البروتينات المهندسة

لقوانين الكيمياء الطبيعية وليس لقيود النشوء. وتسطيع هندسة البروتينات أن تستكشف آفاق ومجالات جديدة.

التطبيقات الاكلينيكية:

هناك اجماع فى الرأى على أن معرفة التتابعات فى الجينوم البشرى وجينوم العديد من الكائنات الأخرى سوف يؤدى الى تحسن كبير فى صحة البشر. وتطبيقات ذلك يمكن أن تتضمن الآتى:

1- تشخيص الأمراض ومخاطرها:

يمكن أن يكشف تتابع الدنا عن غياب جين معين، أو طفرة. يؤدى تعريف تتابعات جين محدد مصاحبا لمرض ما الى سرعة ودقة تشخيص الحالة.

غالبا ما تكون العلاقة بين طبيعة الجين ومخاطر المرض من الصعوبة تحديدها. تعتمد بعض الأمراض مثل الحساسية على تداخلات بين جينات عديدة بالاضافة الى عوامل بيئية. وفى حالات أخرى يكون الجين موجودا وسليما الا أن حدوث طفرة فى مكان ما قد يغير من تعبير ذلك الجين وتوزيعه بين الأنسجة. مثل تلك الأشياء غير السوية يجب اكتشافها بواسطة قياسات نشاط البروتين. كما أن قياس طرز تعبير البروتين تعتبر وسيلة هامة لقياس مدى الاستجابة للعلاج.

2- وراثة الاستجابة الى العلاج المتخصص:

نظرا لاختلاف الناس فى قدرتها على تمثيل الأدوية فإن المرضى المختلفة قد تحتاج الى جرعات مختلفة لعلاج نفس الحالة. يسمح تحليل التتابع اختيار الدواء والجرعة التى تناسب مع كل مريض وذلك فى اطار مجال ينمو بسرعة يعرف باسم Pharmacogenomics. ويصبح فى استطاعة الأطباء تقادى تجريب طرق علاجية مختلفة والتى لها آثار جانبية خطيرة بالاضافة الى التكلفة الباهظة.

٣- تعريف الأماكن المستهدفة للدواء:

المكان المستهدف قد يكون بروتين يتم تعديله وظيفته بالتداخل مع الدواء للتأثير على أعراض المرض ومسبباته. ويلقى التعرف على المكان المستهدف الضوء على الخطوات المتتالية لتصميم الدواء. والأماكن المستهدفة لمعظم الأدوية المستخدمة الآن نصفها مستقبليات وحوالي الربع انزيمات والربع المتبقى هرمونات وحوالي ٧% تعمل على أهداف غير معلومة.

أدى التزايد في ظاهرة مقاومة البكتريا للمضادات الحيوية التي خلق مأساة بخصوص علاج الأمراض. والحاجة الماسة لإيجاد أدوية جديدة مرهونة بقاعدة البيانات اللازمة لإنتاجها. ونتائج دراسة الجينوم قد تؤدي إلى اقتراح أماكن مستهدفة. الاختلافات في الجينوم ومقارنة طرز تعبير البروتين بين سلالات مسببات الأمراض الحساسة والمقاومة للأدوية يمكن أن تحدد البروتين المسئول عن مقاومة الدواء. ويأمل استطاعة دراسة الاختلافات الوراثية بين الخلايا السرطانية والعادية في التعرف على تعبير بروتينات مختلفة تستخدم كأهداف لأدوية مضادة للسرطان.

٤- العلاج بالجينات:

وهو إمكانية استبدال جين - أو على الأقل مد الجسم بنتاجه - محل جين غائب أو مصاب بخلل. كما يتضمن العلاج بالجينات العمل على خفض النشاط الزائد لجين ما. وهناك بعض الحالات المرضية في الإنسان والتي أظهرت نتائج مشجعة بالعلاج الجيني.

ومن التوجهات التي تستخدم لإيقاف تعبير الجين هو ما يطلق عليه "antisense therapy". وهو عبارة عن إنتاج شريط قصير من الدنا أو الرنا الذي يرتبط بأسلوب تتابع متخصص بمنطقة الجين. وهذا الارتباط يتداخل على التوالي مع عمليات النسخ والانتقال. لقد أظهر

هذا النوع من العلاج كفاءة ضد بعض الأمراض مثل cytomegalovirus; and Crohn disease. كما أنه ذهب مباشرة من تتابع المستهدف الى اختصار مراحل عديدة من عمليات تصميم الدواء.

المستقبل The future:

سوف يشهد القرن الحالى تطورا مذهلا فى انتاج وتوزيع وسائل العلاج والعناية بصحة البشر وتتلاشى الحواجز بين قلاع البحث والممارسة الاكلينيكية. ويصبح بإمكان قارئ هذا الكتاب أن يتوصل الى علاج لمرض قاتل، كما يمكن اكتشاف عقاقير ووسائل من شأنها منع حدوث الأورام السرطانية وليس فقط العمل على التحكم فى نموها وانتشارها.

مصادر على الانترنت يمكن الرجوع اليها Web Resources:

Human Genome Project Information :

<http://www.ornl.gov/hgmis/project/info.html>

Genome Statistics:

<http://bioinformatics.weizmann.ac.il/mb/statistics.html>

Taxonomy Sites :

Species : <http://www.sp2000.org>

Tree of life : <http://phylogeny.arizona.edu/tree>

Database of genetics of disease :

<http://www.ncbi.nlm.nih.gov/omim/>

<http://www.geneclinics.org/profiles/all.html>

Lists of databases :

<http://www.infobiogen.fr/services/dbcat/>

<http://www.ebi.ac.uk/biocat/>

List of tools for analysis :

<http://www.ebi.ac.uk./Tools/index.html>

Debate on electronic access to the scientific literature :

<http://www.nature.com/debates/e-access/>

تمرينات محلولة :Solved Problems

Problem 1. For what of following sets of fragment strings does the PERL program mentioned before work correctly ?

(a) Would it correctly recover :

Kate, when France is mine and I am yours, then yours is France and you are mine .

From :

Kate, when France
France is mine
is mine and
and I am\nyours
yours then
Then yours is France
France and you are mine\n

Sample input strings for assembly:

(a) Input data:

Kate, when France
France is mine
is mine and
and I am\nyours
yours then
Then yours is France
France and you are mine\n

(a) Correct answer:

Kate, when France is mine and I am
yours, then yours is France and you are mine.

(b) Would it correctly recover :

One women is fair , yet I am well; another is wise, yet I am well; another virtuous, yet I am well; but till all graces be in one woman, one woman shall not come in my grace .

from :

One woman is
woman is fair,
is fair, yet I am
yet I am well;
I am well; another
another is wise, yet I am well;
yet I am well; another virtuous,
another virtuous, yet I am well;
well; but till all
all graces be

be in one woman,
one woman, one
one woman shall
Shall not come in my grace.

(b) Input data:

One woman is
wo nan is fair,
is fair, yet I am
yet I am well;
I am well; another
another is wise, yet I am well;
yet I am well; another virtuous,
another virtuous, yet I am well;
well; but till all
all graces be
be in one woman,
one woman, one
one woman shall
Shall not come in my grace.

(b) Correct answer:

One woman is fair, yet I am well;
another is wise, yet I am well;
another virtuous, yet I am well;
but till all graces be in one woman,
one woman shall not come in my grace.

(c) Would it correctly recover :

That he is made, 'tis true: 'tis true 'tis pity; And pity 'tis 'tis true.

from :

That he is
is mad, 'tis
'tis true
true: 'tis true 'tis
true 'tis
'tis pity;\n
pity;\nAnd pity
pity 'tis
'tis 'tis
'tis true.\n

(c) Input data:

That he is
is mad, 'tis
'tis true
true: 'tis true 'tis
true 'tis

'tis pity;\n
pity;\nAnd pity
pity 'tis
'tis 'tis
'tis true.\n

(c) Correct answer:

That he is mad, 'tis true: 'tis true 'tis pity;
And pity 'tis 'tis true.

الفصل الثاني

نشوء وتنظيم الجينوم

Genome organization and evolution

الجينوم البكتيري يتأى من جزى دنا منفرد والذى يصل طوله عند فرده الى نحو ٢ ملليمتر (يبلغ قطر الخلية حوالى ٠,٠٠١ ملليمتر). وينظم الدنا فى الكائنات الراقية على كروموسومات حيث تحتوى الخلية العادية فى الانسان على ٢٣ زوجا من الكروموسومات. والكمية الكلية للمعلومات الوراثية لكل خلية - تتابع النيكليوتيدات فى الدنا - تكون ثابتة لدرجة كبيرة لكل أفراد النوع الواحد بينما تختلف بصورة واسعة بين الأنواع المختلفة. وبم أن ليس كل الدنا يشفر الى بروتينات فان كم معلومات تتابع البروتين فى الخلية لا يمكن تقديره بسهولة من حجم الجينوم.

الجينوميك والبروتيوميك Genomics and Proteomics:

البروتيوم Proteomes:

يقدم جينوم الكائن مجموعة كاملة لصفات الفرد. وتعتمد حالة تطور الكائن ونشاطه على المستوى الجزيئى فى أى لحظة على كميات وتوزيع البروتينات. ويعتبر برنامج مشروع البروتيوم من البرامج الضخمة التى تعتنى - باسلوب تكاملى - بطرز تعبير البروتينات فى الأنظمة البيولوجية بشكل يتطابق مع مشروعات الجينوم ويعظمها.

تمييز الجينات فى الجينوم:

تتعرف برامج الكمبيوتر لتحليل الجينوم على أطر قراءة مفتوحة open reading frames (ORFs). ويعرف ORF بأنه منطقة من تتابع الدنا تبدأ

بكود بادئ (ATG) وينتهي بكود ايقاف. وهى منطقة فاعلة لتشفير البروتين.

وعمليات التعريف على تشفير البروتين تتم بواسطة واحد أو الجمع بين اثنان من التوجهات المحتملة التالية:

١- الكشف عن مناطق مشابهة لمناطق تشفير من كائنات أخرى. وتقوم تلك المناطق بتشفير تتابعات حمض أميني مماثل لبروتين معروف أو مماثل علامات تتابعات ناتج الجين Expressed Sequence Tags (ESTs) وبم أن ESTs تشق من الرنا الرسول فانها تناظر جينات معروفة. ومن الضروري أن يتم تحديد التتابع لعدة مئات فقط من قواعد cDNA لتعطي معلومات كافية للتعرف على الجين. ويشابه تحديد الجينات بواسطة ESTs عملية فهرسة القوائد الشعرية والأغاني بالاستعانة بأول بيت فيها.

٢- طرق *Ab initio* والتي تسعى الى التعرف على الجينات من خلال خصائص تتابعات الدنا.

مصادر على الانترنت Web Resources:

HUMAN GENOME INFORMATION

Interactive access to DNA and protein sequences :

<http://www.ensembl.org/>

Images of chromosomes, maps, loci :

<http://www.ncbi.nlm.nih.gov/genome/guide/>

Gene map 99 :

<http://www.ncbi.nlm.nih.gov/genemap99/>

Overview of human genome structure :

<http://hgrep.ims.u-tokyo.ac.jp>

Single-nucleotide polymorphisms :

<http://snp.eshl.org/>

Human genetic disease :

<http://www.ncbi.nlm.nih.gov/Omim/>

<http://www.geneclincs.org/profiles/all.html>

Ethical, legal, social issues :

<http://www.nhgri.nih.gov/ELSI/>

DATABASES OF ALIGNED GENE FAMILIES

Pfam: protein families database :

<http://www.sanger.ac.uk/software/pfam/>

COG: Clusters of orthologous groups :

<http://www.ncbi.nlm.nih.gov/COG/>

HBAGENE : Homologous Bacterial Genes Database :

<http://pbil.univ-lyon1.fr/database/hobacgen.html>

HIOVERGEN: Homologous Vertebrate Genes Database :

<http://pbil.univ-lyon1.fr/databases/hovergen.html>

TAED: The Adaptive Evolution Database :

<http://www.sbc.su.se/~liberles/TAED.html>

GENOME DATABASES

List of completed genomes :

<http://www.ncbi.nlm.nih.gov/PMGifs/Genomes/allorg.html>

<http://www.ebi.ac.uk/genomes/mot/index.html>

<http://pir.georgetown.edu/pirwww/search/genome.html>

Organism-specific databases :

<http://www.unl.edu/stc-95/ResTools/biotools/biotools10.html>

<http://www.-fp.mcs.anl.gov/~gaasterland/genomes.html>

<http://www.lgmp.mrc.ac.uk/GenomeWeb/genome-db.html>

<http://www.bioinformatik.de/cgi->

bin/browse/Catalog/Databases/Genome_projects/

الفصل الثالث

السجلات واسترجاع المعلومات

يتناول هذا الفصل من الكتاب مهارات استرجاع المعلومات والتي تتيح الاستخدام الفعال لبنوك البيانات.

فهرسة قواعد البيانات وخصائص مصطلحات البحث:

الفهرس عبارة عن مجموعة من المؤشرات للمعلومات في قواعد البيانات. عند البحث على شبكة المعلومات الدولية أو من خلال قاعدة بيانات للبيولوجيا الجزيئية يتم ادخال واحد أو أكثر من مصطلحات البحث ويقوم برنامج بالبحث عنها في جداول الفهارس. وهنا يتعرف برنامج الاسترجاع على الموضوعات ذات المحتويات المتعلقة بمجال الاهتمام. على سبيل المثال عند ادخال كلمة حصان horse سوف تحصل على معلومات تختص بموضوعات عديدة تتعلق بالحصان مثل البيولوجيا الجزيئية والتربية والسلالات والقوائد الشعرية حول الحصان وغيرها. وللتركيز على موضوع محدد يسمح نظام استرجاع المعلومات باستخدام كلمات استرشادية. فمثلا بالبحث عن انزيم ديهيدروجينيز للكحول في كبد الحصان horse liver alcohol dehydrogenase سوف يتم التعرف فقط على الموضوعات المحتوية على الكلمات الاسترشادية الأربعة التي تم ادخالها وهي horse و liver و alcohol و dehydrogenase.

لذلك من المهم استخدام قواعد بيانات متخصصة بما في ذلك قواعد بيانات البيولوجيا الجزيئية والتي تفرض تركيب معين على المعلومات لفصل فئات مختلفة من المعلومات.

تحليل البيانات المسترجعة:

أحيانا يكون هناك حاجة الى الحصول على برنامج يسمح باستخدام النتائج - المسترجعة من عمليات بحث في قواعد البيانات - كمدخلات. على سبيل المثال، تستخدم نتائج تتابع بروتين ما كمدخلات للتقريب بواسطة برنامج PSI- BLAST. وفي هذه الحالة يجب تغذية البرنامج بنتيجة التتابع يدويا. الا أنه كما هو الحال في التقريب في قواعد البيانات المتعددة، تتيح نظم المعلومات المسترجعة - في البيولوجيا الجزيئية - تسهيلات لبدء تلك العمليات.

السجلات Archives:

بالرغم من أن معرفتنا ببيانات التتابع والتركيب البيولوجي لم تكتمل بعد، الا أن حجم تلك البيانات ينمو سريعا. ويعمل العديد من العلماء على توفير البيانات أو تنفيذ مشروعات بحثية لتحليل النتائج. كما تقوم هيئات معنية بقواعد البيانات بأرشفة وتوزيع البيانات.

وتجرى عمليات أرشفة البيانات المتعلقة بالمعلوماتية الحيوية بواسطة مجاميع بحثية مهتمة بالعلوم ذات الصلة.

وتتضمن عملية الجمع الأولى للبيانات المتعلقة بالجزيئات البيولوجية الكبيرة مايلي:

- تتابعات الحمض النووي بما في ذلك مشروعات الجينوم.
- تتابعات الحمض الأميني للبروتينات.
- تراكيب البروتين والحمض النووي.
- التراكيب البلورية للجزيئات الصغيرة.
- وظائف البروتين.
- طرز تعبير الجينات.
- المراجع والمؤلفات.

قواعد بيانات تتابع الحمض النووي:

يتكون الأرشيف العالمي لتتابع الحمض النووي من شراكة ثلاثية وهى: المركز القومى لمعلومات التكنولوجيا الحيوية بالولايات المتحدة الأمريكية، ومكتبة البيانات لمعهد اليومعلوماتية الأوروبى بالملكة المتحدة (EMBL)، وبنك اليابان لمعلومات الدنا (المعهد القومى للوراثة باليابان). ويتم تبادل البيانات بين المجموعة يوميا. ولذلك فان البيانات الأولية تكون متطابقة بالرغم من اختلاف نماذج تخزين البيانات وطبيعة تفسيرها. وتقوم قواعد البيانات هذه بمعالجة وأرشفة وتوزيع تتابعات الدنا والرنا المجمعة من مشروعات الجينوم والمنشورات العلمية وتطبيقات براءات الاختراع. وتطلب المجلات العلمية ايداعات جديدة لتتابعات جديدة للنكليوتيدات فى قاعدة البيانات كشرط لنشر المقالة وذلك للتأكد من حرية توفير تلك البيانات الأساسية. وكذلك الحال بالنسبة لتتابعات الحمض الأمينى وتراكيب الحمض النووى والبروتين.

تتكون قواعد بيانات تتابع الحمض النووى من مجموعات من المدخلات. لكل مدخل ملف نص يحتوى على البيانات والتفسيرات لكل تتابع مجلور. العديد من المدخلات تكون مجمعة من عدة أوراق علمية منشورة تقرر أجزاء متداخلة لتتابع كامل.

فيمايلى مثال لمدخلات تتابع دنا من مكتبة بيانات EMBL متضمنة تفسيرات وبيانات تتابع جين مثبط التربسين البنكرياسى:

The EMBL Data Library entry for the Bovine pancreatic trypsin inhibitor gene.

```
ID BTBPTIG standard; DNA; MAM; 3998 BP.
>X
FC X03365; K00966;
>X
EV X03365.1
>X
FT 18-NOV-1986 (Rel. 10, Created)
LT 20-MAY-1992 (Rel. 31, Last updated, Version 3)
XX
DE Bovine pancreatic trypsin inhibitor (BPTI) gene
XX
KW Alu-like repetitive sequence; protease inhibitor; trypsin inhibitor.
XX
OS Bos taurus (cow)
```


OC Eukaryota; Metazoa; Chordata; Craniata; Vertebrata;
 Euteleostomi; Mammalia;
 OC Eutheria; Cetartiodactyla; Ruminantia; Pecora; Bovoidea;
 Bovidae; Bovinae;
 OC Bos.
 XX
 RN [1]
 RP 1-3998
 RX MEDLINE; 86158754.
 RA Kingston I.B., Anderson S.;
 RT "Sequences encoding two trypsin inhibitors occur in strikingly
 similar
 RT genomic environments";
 RL Biochem. J. 233:443-450(1986).
 XX
 RN [2]
 RX MEDLINE; 84070725.
 RA Anderson S., Kingston I.B.;
 RT "Isolation of a genomic clone for bovine pancreatic trypsin
 inhibitor by
 RT using a unique-sequence synthetic dna probe";
 RL Proc. Natl. Acad. Sci. U.S.A. 80:6838-6842(1983).
 XX
 DR SWISS-PROT; P00974; BPT1_BOVIN.
 XX
 CC Data kindly reviewed (08-DEC-1987) by Kingston I.B.
 XX
 FH Key Location/Qualifiers
 FH
 FT source 1..3998
 FT /db_xref="taxon:9913"
 FT /organism="Bos taurus"
 FT misc_feature 795..800
 FT /note="pot. polyA signal"
 FT misc_feature 835..839
 FT /note="pot. polyA signal"
 FT repeat_region 837..847
 FT /note="direct repeat"
 FT misc_feature 930..945
 FT /note="sequence homologous to Alu-like
 consensus seq."
 FT repeat_region 1035..1045
 FT /note="direct repeat"
 FT misc_feature 2456..2461
 FT /note="pot. splice signal"
 FT CDS 2470..2736
 FT /db_xref="SWISS-PROT:P00974"
 FT /note="put. precursor"
 FT /protein_id="CAA27062.1"
 FT
 /translation="PSLFNRDPPIPAAQRPDFCLEPPYTGPKARIIRYFYNAKAGLCQ
 FT TFVYGGCRAKRNNFKSAEDCMRTCGGAIGPWGKTGGRAEGEGKG"
 FT misc_feature 2488..2489
 FT /note="pot. intron/exon splice junction"
 FT misc_feature 2506..2507
 FT /note="pot. intron/exon splice junction"
 FT CDS 2512..2685
 FT /db_xref="SWISS-PROT:P00974"
 FT /note="trypsin inhibitor (aa 1-58)"
 FT /protein_id="CAA27063.1"
 FT
 /translation="RPDFCLEPPYTGPKARIIRYFYNAKAGLCQTFVYGGCRAKRNNF
 FT KSAEDCMRTCGGA"
 FT misc_feature 2698..2699
 FT /note="pot. exon/intron splice junction"
 FT misc_feature 3690..3695
 FT /note="pot. polyA signal"

```

FT misc_feature 3729..3733
FT /note="pot. polyA signal"
XX
SQ Sequence 3998 BP; 1053 A; 902 C; 892 G; 1151 T; 0 other;
aattctgata atgcagagaa ctggaagga gttctgattg ttctgcttga
ttaaatgggt 60
tgtaacagga tagtgtcttg tectgatcct agcattcata tgggtgtgtg
tcggggcaa 120
gtcactcgca gtttcttcaac ctgaacaggg ggaccagggt acatgagttt
cttaaaagat 180
taccagtcac gagtatgaag agtttacct ttctctgatca atgacgtcca
ttcccatca 240
aaatatttta gtccaaaaga ctcatctatc taatgtagat cattttctca
ccaccctct 300
aaaaaattta tctttcagat atgatcattt ctctattatg aaattaatca
gagagttgag 360
tgacagctga gtgtcttccc tccaaggca actgcaggaa gagcaagaaa
tgcaatactt 420
ttctatgagt ttgctctggt ggccaagact gctttttcca ggctggtaca
atagtaatca 480
aatctcaaag atattcttct ttctctctgg ccagactatt attttatttt
cctatcaaga 540
tatagaaagt tagaagtaga ctcataatta tataggcagg cctcatcatc
aaatagacta 600
acaagaattt tattttatct gccttttcaa tgactgtgca ctggcatga
ggatgaaatg 660
ggagatttat tcctttgata aatattcatg aaatacttat gctttttgtc
cctaaaaagc 720
atatttcttg atataggaaa acagctgtaa acaaaaggta gtaaaataat
atgccttcta 780
agagggatac agacaataaa gacgggggag gattcctata ccaggtcag
atgagtgtca 840
tgaggaaggt gagttatggg gttcaggatg ctgtagagga tcagggaaac
cctctgtgat 900
gaggagacat taagcagaag ctgccaaaaa ggagcctggt gtgtttgagc
acagccagg 960
accaggggtg ctggagctga gtgggtgagg ggaggggagt ggaaggggat
geagcagaga 1020
ggccatgggg gcaggtcag aggaacctc taggacttta taaggataaa
aatttgactc 1080
tgagagagct gggaaaccac tgagggactg gtcgttgaa caacgagata
ggactggagt 1140
ttaacaggg tccttgagac tgcagtgtgg agcgtggcct ctagggggcg
aaaagcaggg 1200
acagggggcc cctgggcagg tggctgcagg ggtccagtga gctatgatgg
agaatatata 1260
cctgtgttgt tcctgggtt gattccagtt ctcttgaata acctgatga
cttgctatat 1320
taagatatct ggggaggctt catcacaat catgattcat taaatcttta
gfcatttctg 1380
attgattcaa cctccaatcc ctctcccctt cctgagatgt ggatacaatg
aaagaagtag 1440
gaatgaaaat tcccacacc aactcaggca gttgtttccc ctgacaactt
atccccattc 1500
ttggctttgc ttgaggcttt acaaaactca tctccctcac atgataaagg
actccccctt 1560
gctctcatct cttaggaaat tccaatggtt taggagctct gtggcaggaa
tgggatgcag 1620
accaagttaa tatttctttt ataagtcaca gtatcaatat ttctcaata
ttctattatt 1680
ccagtctcca tgaggaacc aaagtaaca ccggtgtgtt ttctaccatg
cctttctcca 1740
tttatggcat gatttctca cacttttgta atagtgcgg gtcacgcagg
cctatcaacc 1800
attggctggc atccaggtgg gcacctcat caagggataa ctgtaaatga
gcaacccttg 1860

```

gtggccagt	tcagccattg	ccactgttgt	agccacagtg	ggctcttcgc		
cctcccgttc	1920	ttttgtataa	aaggaacagg	aatttatact	gtgggaagat	ggttttcttg
agacagtagc	1980	atgctatcat	ctccgtgggc	ccaatttcca	attaaaaatg	ttattcctag
ttccagcaac	2040	tcttctccgg	attattggct	ggccctgagg	tgagcagaat	gagactgggc
tcagtgatgc	2100	tttcttaacg	gtggaagttt	ccaccacaca	catacataga	aagcatagta
ttaaaaaagc	2160	cgtggatcca	ttgtccagct	ccagtaattt	ctatacatgg	agagtatttt
tatatgtgtc	2220	cctcttttgt	gttacttttg	aaactcatca	gtagcatcat	gctaattaat
gcataaacat	2280	tcataaatgg	catgtaatta	tttataatat	tgccctgtca	ttgtcacacc
taacaacatt	2340	aataataatg	tcctggaaag	cagggtgtca	aaaggccttt	tcacgtttca
cacttctgcc	2400	caccccccat	cactctctat	cacaaactgg	tggttttagt	tyttcatctt
gtagactgag	2460	ctgtgatgac	cttccctctt	taaccgagat	cctcccacc	ctgcagccca
ggggcctgac	2520	ttctgcttag	agcctccata	taagggtccc	tgeaaggcca	gaattatcag
atacttctac	2580	aacgccaagg	ctgggctctg	ccagaccttt	gtatatggcg	gctgcagagc
taaaagaaac	2640	aatctcaaga	gcgcagagga	ctgcatgagg	acctgtggtg	gtgctattgg
gcccctgggt	2700	aagacagggg	gcagggcaga	gggagagggg	aagggtagg	gaaagtgggt
gcgctcagaa	2760	ggccacacac	ctttccaaaa	aagtgatatt	tttcccttgt	tgcctcccaa
gagaagtgtc	2820	agaagtatcc	gtggattgag	catgtcctcc	atggaccagc	ttggtgaaag
gccaccccc	2880	agaagccctg	tcatacataat	ctgagcctac	tcacatgctc	ccatttttca
gatgggaaca	2940	ctgagtcagt	cactctgcag	agcaagtctg	gagtgtcctc	cagtccccca
cctcagcctg	3000	gaaaactccc	ttgtttattg	ttggttatcc	tggtcctggg	aggactgtgg
ttgcgcattt	3060	ctgggatggt	ctaggacctg	tcagggtgga	cagtgtccag	gtctgggcct
tcagagatgt	3120	cattcagcaa	gttccctttct	ttttacagag	aacctgtgaa	ctgtgtctcc
ctgagatgct	3180	gaagtatgag	gaggaccacc	ccaaggctgg	cctctatctg	cttctgaaaa
atttcagcct	3240	ccttttattt	cttctcaacc	ctccccctct	cagcagaaat	ctgtctcttt
ccctctctca	3300	caggctccact	tactttagec	ctatctcacc	cagtttgtct	taagcaccat
gaaagcaaat	3360	cttccctttg	tcccctcacac	ttcccacaat	ttctggcaca	aaggagaagg
tcagaaata	3420	ttggaggaag	gaaggaatga	agttcccact	gactggagca	tctgtagagt
ctjagattta	3480	aatctggatt	ctgtctctaa	tcttctctct	cacggcatcc	ttaccttcat
cctccacccc	3540	accatcactg	ctctccctct	actggcgaaa	gtagaatttc	catcatcgag
ttttcagctc	3600	agtggtgggg	gaggtctttt	catgaacgaa	acctctcct	cacattgatt
tgaaggtctg	3660	tggtctcaaa	gagtctggcc	ttatctttaa	ataaattcat	attttaatta
aactaactgg	3720	agtggattgt	gttgtttgca	actaagaacc	ttaaccata	ggttccatgg
aaacgggtgt	3780	ctttctcatt	ttatgcagat	gggtgggcag	ctctccatca	cctctctcca
gactcagccc	3840					

```
taccaagtag   aaggagccaa   cccttacac   tgacatctac   ctcttatggc
cgtgccagtg   3900
tacaatgaaaa  actggatgag   agacacctca  acaagaaaac   ttttgcctt
cacttcttgg   3960
gccaggtcaa  actttgggggt  gtgttatttc  cctgaatt     3998
//
```

السطور التي تبدأ ب FT هي عبارة عن مكون لتفسير المدخل الذي يقرر صفات مناطق متخصصة (تتابعات كودية CDS) بحيث تقرأ بواسطة برامج كمبيوتر. على سبيل المثال: لترجمة كود منطقة الى تتابع حمض أميني، توجد نماذج متحكم فيها بعناية ومفردات مقيدة. ويكون من المهم وجود مفردات لغوية محكمة وقواميس وكلمات استرشادية وجداول مميزة حتى يمكن انشاء روابط بين قواعد بيانات مختلفة.

قواعد بيانات الجينوم:

بالرغم من أن تتابعات الجينوم تشكل مدخلات الأرشيفات المرجعية لتتبع الحمض النووي، فإن العديد من الأنواع لها قواعد بيانات خاصة والتي تجمع تتابع الجينوم وتفسيره مع باقي البيانات ذات العلاقة بالنوع.

مصادر على الانترنت Web Resources لروابط قواعد بيانات للكائنات:

<http://www.unl.edu/stc-95/ResTools/biotools/biotools10.html>

<http://www.-fp.mcs.anal.gov/~gaasterland/genomes.html>

<http://www.hgmp.mrc.ac.uk/GenomeWeb/genome-db.html>

http://www.bioinformatik.de/cgi-bin/browse/Catalog/Database/Genome_projects/

قواعد بيانات تتابع البروتين:

تتأى بيانات تتابع الحمض الأميني من ترجمة تتابعات الحمض النووي. يتعاون المعهد السويسري للبيومعلوماتية مع مكتبة بيانات EMBL لتزود قاعدة بيانات تفسيرية لتتابعات الحمض الأميني تسمى SWISS-PROT. كما توجد قاعدة بيانات لتتابع البروتين تنتجها The PIR International والتي تشكل مجموعات في المؤسسة القومية للبحوث الطبية بجامعة جورج تاون واشنطن - الولايات المتحدة الأمريكية، ومركز معلومات

ميونيخ لتتابعات البروتين (MPIS) بألمانيا، وقاعدة بيانات اليابان الدولية
لمعلومات البروتين.

وفيمائلي نموذج لمدخلات تتابع الحمض الأميني لبروتين لمثبط التربسين
البنكرياسي في قاعدة بيانات PIR:

PIR entry for the amino-acid sequence of Bovine pancreatic trypsin inhibitor

```
ENTRY          TIBO #type complete iProClass View of TIBO
TITLE          basic proteinase inhibitor precursor - bovine
ALTERNATE_NAMES aprotinin; basic pancreatic trypsin inhibitor;
BPTI;
                cationic kallikrein inhibitor; inhibitor IV
ORGANISM       #formal_name Bos primigenius taurus #common_name
cattle
                #cross-references taxon:9913
DATE           24-Apr-1984 #sequence_revision 22-Jul-1994
#text_change
                16-Jun-2000
ACCESSIONS     S00277; A30333; S10546; S02486; S28197; A90162;
A92023;
                A90736; A90927; A34658; A93977; S10062; A01205
REFERENCE      S00274
#authors       Creighton, T.E.; Charles, I.G.
#journal       J. Mol. Biol. (1987) 194:11-22
#title        Sequences of the genes and polypeptide precursors
for two
                bovine protease inhibitors.
#cross-references MUID:87283904
#accession     S00277
                ##molecule_type DNA; mRNA
                ##residues 1-100 ##label CR2
                ##cross-references GB:M20934; GB:X05274; NID:g162767;
                PIDN:AAD13685.1; PID:g162769
REFERENCE      A90926
#authors       Creighton, T.E.; Charles, I.G.
#journal       Cold Spring Harb. Symp. Quant. Biol. (1987) 52:511-
519
#title        Biosynthesis, processing, and evolution of bovine
pancreatic trypsin inhibitor.
#cross-references MUID:88295740
#accession     A30333
                ##molecule_type DNA
                ##residues 1-100 ##label CRE
                ##cross-references GB:M20934; GB:X05274; NID:g162767;
                PIDN:AAD13685.1; PID:g162769
REFERENCE      S10546
#authors       Kingston, I.B.; Anderson, S.
#journal       Biochem. J. (1986) 233:443-450
#title        Sequences encoding two trypsin inhibitors occur in
strikingly similar genomic environments.
#cross-references MUID:86158754
#accession     S10546
                ##molecule_type DNA
                ##residues 34-97 ##label KIN
REFERENCE      S02485
#authors       Fioretti, E.; Angeletti, M.; Fiorucci, L.; Barra,
D.;
                Bossa, F.; Ascoli, F.
#journal       Biol. Chem. Hoppe-Seyler (1988) 369(Suppl.):37-42
```

```

#title      Aprotinin-like isoinhibitors in bovine organs.
#cross-references MUID:89076531
#accession  S02486
  ##molecule_type protein
  ##residues 36-93 ##label FIO
REFERENCE  S28197
#authors   Ikekita, M.; Jone, C.S.; Kamo, M.; Tsugita, A.;
Kizuki,
           K.; Moriya, H.
#journal   Protein Seq. Data Anal. (1992) 5:7-11
#title     Purification and characterization of the major
cationic
           kallikrein inhibitor in bovine pituitary gland.
#cross-references MUID:93150003
#accession  S28197
  ##molecule_type protein
  ##residues 36-93 ##label IKE
REFERENCE  A90162
#authors   Kassell, B.; Laskowski, M.
#journal   Biochem. Biophys. Res. Commun. (1965) 20:463-468
#title     The basic trypsin inhibitor of bovine pancreas. V.
The
           disulfide linkages.
#cross-references MUID:66083012
#contents  annotation; disulfide bonds
#accession  A90162
  ##molecule_type protein
  ##residues 36-93 ##label KAS
REFERENCE  A92023
#authors   Anderer, F.A.; Hornle, S.
#journal   J. Biol. Chem. (1966) 241:1568-1572
#title     The disulfide linkages in kallikrein inactivator of
bovine
           lung.
#cross-references MUID:66171231
#contents  annotation; disulfide bonds
#accession  A92023
  ##molecule_type protein
  ##residues 36-93 ##label AN2
REFERENCE  A90736
#authors   Chauvet, J.; Acher, R.
#journal   Bull. Soc. Chim. Biol. (1967) 49:985-1000
#title     La structure covalente d'un inhibiteur
polypeptidique de
           la trypsine (inhibiteur de Kunitz et Northrop).
#cross-references MUID:68012003
#contents  annotation; disulfide bonds
#accession  A90736
  ##molecule_type protein
  ##residues 36-93 ##label CHA
REFERENCE  A90927
#authors   Dlouha, V.; Pospisilova, D.; Meloun, B.; Sorm, F.
#journal   Collect. Czech. Chem. Commun. (1968) 33:1363-1365
#title     Sequence of residues 18-20 in pancreatic trypsin
inhibitor.
#accession  A90927
  ##molecule_type protein
  ##residues 36-93 ##label DLO
REFERENCE  A93410
#authors   Huber, R.; Kukla, D.; Ruhlmann, A.; Epp, O.;
Formanek, H.
#journal   Naturwissenschaften (1970) 57:389-392
#title     The basic trypsin inhibitor of bovine pancreas. I.
Structure analysis and conformation of the
polypeptide
           chain.
#cross-references MUID:70255230

```

```

#contents      annotation; X-ray crystallography of basic protease
                inhibitor, 2.5 angstroms
REFERENCE      A34658
#authors       Lewis, R.V.; Ray, P.; Coguill, R.; Kruggel, W.
#journal       Biochem. Biophys. Res. Commun. (1990) 167:543-547
#title         Presence of pancreatic trypsin inhibitor in adrenal
                medullary chromaffin cells.
#cross-references MUID:90211226
#accession     A34658
                ##molecule_type protein
                ##residues 36-53,55-81.##label LEW
REFERENCE      A93977
#authors       Anderson, S.; Kingston, I.B.
#journal       Proc. Natl. Acad. Sci. U.S.A. (1983) 80:6838-6842
#title         Isolation of a genomic clone for bovine pancreatic
                trypsin
                inhibitor by using a unique-sequence synthetic DNA
                probe.
#cross-references MUID:84070725
#accession     A93977
                ##molecule_type DNA
                ##residues 'PSLFNRDPPIPA',34-97,'GKTGGRAEGEGKG' ##label AND
                ##cross-references GB:X03365; GB:K00966; NID:g142;
                PIDN:CAA27062.1;
                PID:g1364183
REFERENCE      S00371
#authors       Siekmann, J.; Wenzel, H.R.; Schroeder, W.;
                Tschesche, H.
#journal       Biol. Chem. Hoppe-Seyler (1988) 369:157-163
#title         Characterization and sequence determination of six
                aprotinin homologues from bovine lungs.
#cross-references MUID:88221840
#accession     S10062
                ##molecule_type protein
                ##residues 36-66,'P',68-82,'S',84-93 ##label SIE
                ##experimental_source lung
                ##note the authors designated this protein as isoaprotinin 2
COMMENT        Basic proteinase inhibitor is an intracellular
                polypeptide
                found in many tissues, probably located in granules
                of
                connective tissue mast cells.
GENETICS
#introns       34/1; 98/1
CLASSIFICATION #superfamily basic proteinase inhibitor; animal
                Kunitz-type proteinase inhibitor homology
KEYWORDS       serine proteinase inhibitor
FEATURE
1-20           #domain signal sequence #status predicted
#label
                SIG\
21-35         #domain propeptide #status predicted
#label PRO\
36-100        #product basic proteinase inhibitor
#status
                experimental #label MAT\
40-90         #domain animal Kunitz-type proteinase
inhibitor
                homology #label BPI\
40-90,49-73,65-86 #disulfide_bonds #status experimental\
50            #inhibitory_site Lys (trypsin,
                chymotrypsin,
                kallikrein, plasmin) #status experimental
SUMMARY        #length 100 #molecular_weight 10903

SEQUENCE

```

```

      5      10      15      20      25      30
1 M K M S R L C L S V A L L V L L G T L A A S T P G C D T S N
31 Q A K A Q R P D F C L E P P Y T G P C K A R I I R Y F Y N A
61 K A G L C Q T F V Y G G C R A K R N N F K S A E D C M R T C
91 G G A I G P W E N L
    
```

in the PIR1 section of the Protein Sequence Database, release
71.00,
31-Dec-2001, assembled and annotated by the PIR-International.
Copyright 2000 PIR-International.

PDB structures most related to TIBO:

```

1CBWD (36-93) 100.0%; 1BZ5E (36-93) 100.0%; 9PTI (36-91)
100.0%
1BZXI (36-93) 100.0%; 1B0CB (36-93) 100.0%; 1CBWI (36-93)
100.0%
1B0CD (36-93) 100.0%; 1B0CE (36-93) 100.0%; 6PTI (36-93)
100.0%
1BHCA (36-93) 100.0%; 1BHCC (36-93) 100.0%; 1MTND (36-93)
100.0%
1BHCE (36-93) 100.0%; 1MTNH (36-93) 100.0%; 1BHCG (36-93)
100.0%
1PIT (36-93) 100.0%; 1BHCI (36-93) 100.0%; 1TPAI (36-93)
100.0%
1BPI (36-93) 100.0%; 2HEXA (36-93) 100.0%; 1BTHQ (36-93)
100.0%
2HEXB (36-93) 100.0%; 1BZ5B (36-93) 100.0%; 2HEXC (36-93)
100.0%
1BZ5D (36-93) 100.0%; 2HEXD (36-93) 100.0%; 1B0CC (36-93)
100.0%
2HEXE (36-93) 100.0%; 1BHCD (36-93) 100.0%; 2KAI (36-93)
100.0%
1BHCH (36-93) 100.0%; 2PTCI (36-93) 100.0%; 1BTHP (36-93)
100.0%
2TGPI (36-93) 100.0%; 1BZ5C (36-93) 100.0%; 2TPII (36-93)
100.0%
1BHCB (36-93) 100.0%; 3TGII (36-100) 100.0%; 1BHCJ (36-93)
100.0%
3TGJI (36-100) 100.0%; 1B0CA (36-93) 100.0%; 3TPII (36-93)
100.0%
1BZ5A (36-93) 100.0%; 4PTI (36-93) 100.0%; 1BHCF (36-93)
100.0%
5PTI (36-93) 100.0%; 1FAN (36-93) 98.3%; 1BPT (36-93) 98.3%
1NAG (36-93) 98.3%; 1BTI (36-93) 98.3%; 4TPII (36-93) 98.3%
8PTI (36-93) 98.3%; 1AALA (36-93) 96.6%; 1AALB (36-93) 96.6%
7PTI (36-93) 96.6%; 1BRBI (36-93) 94.8%; 1QLQA (36-93) 93.1%
    
```

SCOP: 1CBW ; 1BZ5 ; 9PTI ; 1BZX ; 1B0C ; 6PTI ; 1BHC ; 1MTN ;
1PIT ; 1TPA
; 1BPI ; 2HEX ; 1BTH ; 2KAI ; 2PTC ; 2TGP ; 2TPI ; 3TGI ; 3TGJ ;
3TPI ; 4PTI
; 5PTI ; 1FAN ; 1BPT ; 1NAG ; 1BTI ; 4TPI ; 8PTI ; 1AAL ; 7PTI ;
1BRB ; 1QLQ

CATH: 1CBW ; 1BZ5 ; 9PTI ; 1BZX ; 1B0C ; 6PTI ; 1BHC ; 1MTN ;
1PIT ; 1TPA
; 1BPI ; 2HEX ; 1BTH ; 2KAI ; 2PTC ; 2TGP ; 2TPI ; 3TGI ; 3TGJ ;
3TPI ; 4PTI
; 5PTI ; 1FAN ; 1BPT ; 1NAG ; 1BTI ; 4TPI ; 8PTI ; 1AAL ; 7PTI ;
1BRB ; 1QLQ

FSSP: 1CBW ; 1BZ5 ; 9PTI ; 1BZX ; 1B0C ; 6PTI ; 1BHC ; 1MTN ;
 1PIT ; 1TPA
 ; 1BPI ; 2HEX ; 1BTH ; 2KAI ; 2PTC ; 2TGP ; 2TPI ; 3TGI ; 3TGJ ;
 3TPI ; 4PTI
 ; 5PTI ; 1FAN ; 1BPT ; 1NAG ; 1BTI ; 4TPI ; 8PTI ; 1AAL ; 7PTI ;
 1BRB ; 1QLQ

MMDB: 1CBW ; 1BZ5 ; 9PTI ; 1BZX ; 1B0C ; 6PTI ; 1BHC ; 1MTN ;
 1PIT ; 1TPA
 ; 1BPI ; 2HEX ; 1BTH ; 2KAI ; 2PTC ; 2TGP ; 2TPI ; 3TGI ; 3TGJ ;
 3TPI ; 4PTI
 ; 5PTI ; 1FAN ; 1BPT ; 1NAG ; 1BTI ; 4TPI ; 8PTI ; 1AAL ; 7PTI ;
 1BRB ; 1QLQ

ALIGNMENTS containing TIBO:

FA2061 basic proteinase inhibitor - 328.8 1.0
 SA0572 basic proteinase inhibitor superfamily 328.8
 M01603 basic proteinase inhibitor - 1561.0 1.0

Associated Alignments:

DA1053 animal Kunitz-type proteinase inhibitor homology

يوجد اثنان من قواعد البيانات المشاركة ل SWISS-PROT وهما
 .ENZYME DB & PROSITE

قواعد البيانات المشاركة ل SWISS-PROT

يقوم ENZYME DB بتخزين المعلومات التالية عن الانزيمات:

- EC Number وهو وسيلة تعريف رقمية تم اقرارها بواسطة مفوضية
 الانزيمات المنبثقة من الاتحاد الدولي للكيمياء الحيوية والبيولوجيا
 الجزيئية، (يمكن الرجوع الى الموقع التالي:
<http://www.chem.qmw.ac.uk/iubmb/enzyme/>)
- الاسم الموصى به.
- الاسم البديل، ان وجد.
- النشاط المحفز.
- المرافقات.
- مؤشرات الى SWISS-PROT وبنوك المعلومات الأخرى.
- مؤشرات الى الأمراض المصاحبة للخلل في النشاط الانزيمي.

والمثال التالي يوضح المخلات في ENZYME DB حيث:

التعريف	=	ID
الوصف = الاسم الرسمي	=	DE
الاسم البديل	=	AN
النشاط التحفيزي	=	CA
المرافقات	=	CF
التعليقات	=	CC
مرجع قاعدة البيانات (SWISS-PROT)	=	DR

A Sample Entry in ENZYME DB

```
ID 1.14.17.3
DE PEPTIDYLGLYCINE MONOOXYGENASE.
AN PEPTIDYL ALPHA-AMIDATING ENZYME.
AN PEPTIDYLGLYCINE 2-HYDROXYLASE.
CA PEPTIDYLGLYCINE + ASCORBATE + O(2) = PEPTIDYL(2-HYDROXYGLYCINE)+
CA DEHYDROASCORBATE + H(2)O.
CF COPPER.
CC -!- PEPTIDYLGLYCINES WITH A NEUTRAL AMINO ACID RESIDUE IN THE
CC PENULTIMATE POSITION ARE THE BEST SUBSTRATES FOR THE ENZYME.
CC -!- THE ENZYME ALSO CATALYZES THE DISMUTATION OF THE PRODUCT TO
CC GLYOXYLATE AND THE CORRESPONDING DESGLYCINE PEPTIDE AMIDE.
DR P10731, AMD_BOVIN ; P19021, AMD_HUMAN ; P14925, AMD_RAT ;
DR P08478, AMD1_XENLA; P12890, AMD2_XENLA;
```

: PROSITE

تحتوى على طرز لمتبقيات مجاميع البروتينات مثل الطراز (البصمة أو القالب) الذى يظهر عادة فى عائلة من البروتينات القريبة بسبب متطلبات أماكن الارتباط والتي تحد من التحول فى عائلة بروتين ما.
PIR وقواعد البيانات المشاركة:

يدير PIR عدة قواعد بيانات خاصة بالبروتينات مثل:

- PIR-PSD : وهى قاعدة البيانات الأساسية لتتابع البروتين.
- iProClass : تقسيم البروتينات طبقا للتركيب والوظيفة.

- ASDB : قاعدة بيانات التفسير والمماثلة، وكل مدخل متصل بقائمة من التتابعات المماثلة.
 - P/R-NREF : تجميع لأكثر من ٨٠٠٠٠٠٠٠ تتابع بروتين من مصادر متاحة.
 - NRL3D : قاعدة بيانات لتتابعات وتفسيرات لبروتينات معروفة التركيب وموجودة ببنك معلومات البروتين.
 - ALN : قاعدة بيانات لمحازاة تتابع البروتين.
 - RESID : قاعدة بيانات لتحورات التركيب التساهمي للبروتين.
- وقد أنشأ PIR موقع IESA (Integrated Environment for Sequence Analysis) لاسترجاع المعلومات والمعاملة الحسابية لها. ويمكن الرجوع الى موقع PIR على الوب:

<http://pir.georgetown.edu>

قواعد بيانات التراكيب:

يقوم أرشيف قواعد بيانات التركيب بتفسير وتوزيع مجاميع من المنسقات الذرية. ويعتبر بنك بيانات البروتين (PDB) Protein Data Bank من أفضل قواعد البيانات الخاصة بتراكيب الجزيئات البيولوجية الكبيرة. ويحتوى البنك على تراكيب بروتينات وأحماض نووية و قليل من الكربوهيدرات. ويدار البنك بواسطة the Research Collaboratory for Structural Bioinformatics (RCSB).

والموقع الإلكتروني الأساسي لبنك بيانات البروتين هو: <http://www.rcsb.org>. وهناك مواقع رسمية للبنك لأماكن عديدة حول العالم منها أوروبا، سينغافورة، اليابان، والبرازيل.

وتحتوى الصفحة الأساسية للبنك على روابط اتصال ملفات البيانات ذاتها، ومواد تعليمية، أخبار، نشرات دورية وبرامج بحث متخصصة لمعالجة التراكيب.

مثال: يمكن الحصول على بيانات بروتين ثيورودوكسين في بكتريا *E. coli* (Protein data bank entry 2trx, *E. coli* thioredoxin) من خلال الموقع التالي:

<http://www.rcsb.org/pdb/cgi/explore.cgi?job=download;pdbid=2TRX;page=0;pid=197301010382789&opt=show&format=PDB>

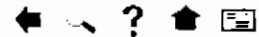
وتضمن المدخلات البيانات التالية:

- نوع البروتين ومصدره.
- مراجع لوصف طرق تحديد التركيب.
- بيانات تفصيلية لتجارب تحديد التركيب متضمنة نتائج الأشعة السينية واحصائيات للشكل الفراغي للتركيب.
- تتابع الحمض الأميني.
- الجزيئات الاضافية مثل المرافقات، المثبطات، جزيئات الماء.
- التركيب الثانوي: الشكل المطوى، المنفرد.
- المواقع ثنائية الكبريت.
- المنسقات الذرية.

وقد حدد بنك بيانات البروتين كود تعرف رباعي لكل تركيب موجود بالبنك، يبدأ برقم (من 1 الى 9). وأصبح من السهل استرجاع تركيب ما وذلك باستخدام كود التعريف الخاص به. حيث يمكن ادخال الكود على الصفحة الرئيسية لموقع RCSB ثم اختيار Explorer للحصول على ملخص لبيانات التركيب من صفحة واحدة كما هو مبين:

Protein data bank entry Structure Explorer - 2trx, *E. coli* thioredoxin

<http://www.rcsb.org/pdb/cgi/explore.cgi?pid=197301010382789&page=0&pdbid=2TRX>



Summary Information

Summary Information

[View Structure](#)

[Download/Display File](#)

[Structural Neighbors](#)

[Geometry](#)

[Other Sources](#)

[Sequence Details](#)

[Crystallization Info](#)

Title: Crystal structure of thioredoxin from Escherichia coli at 1.68 Å resolution.

Compound: Thioredoxin

Authors: S. K. Katti, D. M. LeMaster, H. Eklund

Exp. Method: X-ray Diffraction

Classification: Electron Transport

Source: Escherichia coli

Primary Citation: [Katti, S. K., LeMaster, D. M., Eklund, H.:](#) Crystal structure of thioredoxin from Escherichia coli at 1.68 Å resolution. *J Mol Biol* 212 pp. 167 (1990)

Deposition Date: 19-Mar-1990 **Release Date:** 15-Oct-1991

Resolution [Å]: 1.68 **R-Value:** 0.165

Space Group: C 2

Unit Cell: dim [Å]:

<i>a</i>	89.50	<i>b</i>	51.06	<i>c</i>	60.45
<i>alpha</i>	90.00	<i>Beta</i>	113.50	<i>gamma</i>	90.00

Polymer Chains: A, B **Residues:** 216

Atoms: 1842

Chemical Components:
("HET" groups)

ID (needs Rasmol)	Name	Formula	Retrieve All PDB IDs Containing
CU	COPPER (II) ION	2(Cu)	CU
MPD	2-METHYL-2,4-PENTANEDIOL	7(C ₆ H ₁₄ O ₂)	MPD

CATH: [Structural Classification](#)

PDBSum: [Summary of PDB Structure](#)

SCOP: [Structural Classification](#)

[RCSB](#)

البيومعلوماتية

٦٢

- والصفحة السابقة بها روابط اتصال للمعلومات التالية:
- مراجع للبيانات الموجودة من خلال قاعدة بيانات المراجع PubMed.
 - صور للتراكيب ويحتاج ذلك وجود برنامج لظهار الصور.
 - مداخل لملفات البيانات.
 - قوائم بالتراكيب الأخرى ذات الصلة طبقاً لتصنيف تراكيب البروتين.
 - التحليل الكيميائي الفراغي، وتوزيع أطوال الروابط والزوايا.
 - مصادر لمعلومات حول المدخلات.
 - التتابع والتراكيب الثانوية.
 - تفاصيل حول الشكل البللوري وطرق إنتاج البللور.
- وفي حالة عدم معرفة كود التعرف يمكن استخدام أداة بسيطة من على الصفحة الرئيسية لبنك بيانات البروتين وهي Search Lite والتي تسمح باستخدام كلمات استرشادية key words.
- مصادر على الانترنت Web Resources لتراكيب البروتين والحمض النووي:

الصفحة الرئيسية لبنك بيانات البروتين:

<http://www.rcsb.org>

الصفحة الرئيسية لقاعدة بيانات EBI التركيب الجزيئي الماكرو:

<http://msd.ebi.ac.uk/>

الصفحة الرئيسية لـ BioMagResBank

<http://www.bmrb.wisc.edu/>

للبحث في بنك بيانات البروتين:

الصفحة الرئيسية لـ SCOP (Structural Classification of Protein)

<http://scop.mrc-lmb.cam.ac.uk/scop/>

قائمة بأماكن التصفح:

http://pdp-browsers.ebi.ac.uk/browse_it.shtml

أداة البحث OCA:

<http://oca.ebi.ac.uk/oca-bin/ocamain>

قاعدة بيانات التركيب الرباعي للبروتين:

<http://pqs.ebi.ac.uk/>

تقارير جودة التركيب:

<http://www.cmbi.kun.nl/gv/pdbreport>

مصادر على الانترنت Web Resources لقواعد بيانات لعائلات معينة
من بروتين:

Protein kinase :

<http://www.sdsc.edu/kinases/>

HIV proteases

<http://www-fbnc.ncifcrf.gov/HIVdb/>

Icosahedral viruses :

<http://mmtsb.scripps.edu/viper/main.html>

قاعدة بيانات علم المناعة Immunology :

IGMT : (International ImMunoGene Tics database):

<http://imgt.cines.fr>

وهي قاعدة بيانات متكاملة عالية الجودة، ومتخصصة في:

و Immunoglobulins (Ig)، و T-Cell receptors (TcR) ، و

Histocompatibility Complex (MHC) molecules of all vertebrate species.

KABAT : <http://immuno.bme.nwu.edu./>

وهي قاعدة بيانات لتتابعات البروتينات ذات الأهمية المناعية.
MHCPEP : <http://wehih.wehi.edu.au/mhcpep/>
وهي قاعدة بيانات لمعقد البيبتيدات الرئيسية المرتبطة والمتوافقة
هستولوجيا.

مجموعات من روابط الاتصال لقواعد بيانات لعائلات معينة من
البروتين:

<http://www2.ebi.ac.uk/msd/Links?family.shtml>

قواعد بيانات ناتج تعبير الجين والبروتيومكس :proteomics databases

من المعروف أن الدنا يصنع الرنا ويقوم الرنا بصنع البروتين. وتحتوي
قواعد بيانات الجينوم على تتابعات الدنا. كما تسجل قواعد بيانات
التعبير الجيني مقاييس مستويات الرنا الرسول mRNA وذلك عبر short
(terminal sequences of cDNA synthesized from mRNA) ESTs
وصف طرز نسخ الجين. بينما تسجل قواعد بيانات البروتيومات
proteomics مقاييس للبروتينات مع وصف طرز ترجمة الجين.

تقدم مقارنات طرز التعبير الجيني معلومات موثقة عن:

- (١) وظيفة وآلية فعل نواتج الجين.
- (٢) كيف تنسق الكائنات عملية التحكم في العمليات الأيضية تحت
ظروف مختلفة. على سبيل المثال: الخميرة تحت ظروف هوائية
ولا هوائية.
- (٣) الاختلافات في تعبئة الجين أثناء مراحل مختلفة من دورة الخلية
وتطور الكائن.
- (٤) آليات المقاومة للمضادات الحيوية في البكتيريا وبالتالي اقتراح
أماكن مستهدفة لاستحداث أدوية.
- (٥) الاستجابة لتحدي طفيل ما.

(٦) الاستجابة لأنواع وجرعات مختلفة من الأدوية للوصول الى علاج فعال.

وهناك عدة قواعد بيانات ل ESTs، تحتوي مدخلات معظمها على مجالات تشير الى موقع نسيج و/ أو مكون تحت خلوى، حالة التطور، ظروف النمو، والمستوى الكمي للتعبير. حاليا تحتوي dbEST من خلال بنك الجين على ما يقرب من تسعة ملايين مدخل من ٣٤٨ نوع.

بعض مجموعات EST متخصصة في أنسجة معينة (عضلات، أسنان) أو أنواع. كما توجد جهود لربط طرز التعبير بمعلومات أخرى حول الكائن. فعلى سبيل المثال ينسق مشروع تطور الجرذان بين بيانات تعبير الجين والتشريح التطوري.

تتيح العديد من قواعد البيانات وسائل للربط بين ESTs في أنواع مختلفة مثل ربط التماثل في الانسان والجرذان، أو العلاقات بين جينات أمراض الانسان وبروتينات الخميرة. وهناك مجموعات أخرى متخصصة في نوع من البروتين مثل سيتوكينات. كما يوجد جهدا هائلا للتركيز على السرطان من حيث تكامل المعلومات عن الطفرات واعادة التنظيم للكروموسومات والتغيرات في طرز التعبير وذلك للتعرف على التغيرات الوراثية أثناء تكوين الأورام ونموها.

بالرغم من العلاقة الشديدة بين طرز النسخ وطرز الترجمة، الا أن القياسات المباشرة للمحتوى البروتيني للخلايا والأنسجة - البروتيوميكس proteomics تتيح معلومات إضافية قيمة. وبسبب عدلات النسخ المختلفة لمختلف الرنا الرسول فان قياسات البروتينات تعطى مباشرة وصفا أكثر دقة لطرز تعبير الجين بالمقارنة بقياسات النسخ. ويمكن اكتشاف تحورات ما بعد الترجمة فقط عن طريق فحص البروتينات.

يتضمن تحليل البروتيوم الفصل والتعريف والتقدير الكمي لبروتينات العينة. ويعتمد ذلك على فرد البروتين بواسطة الجيل ثنائي الأبعاد

وتعريف كل مكون بواسطة مطياف الكتلة. وتقوم قواعد بيانات البروتيوم بتخزين صور الجيل وتفسيراتها على هيئة طرز بروتين. كما تظهر بعض قواعد البيانات صوراً وتسمح باختيار تفاعلي للنقط. وباختيار نقطة معينة تفتح نافذة لمداخل البيانات المتطابقة. ويوجد مدخل بيانات لكل بروتين يسجل المعلومات التالية:

- تعريف البروتين.
- الكمية.
- الوظيفة.
- آلية الفعل.
- طراز التعبير.
- مكان التواجد تحت الخلوى.
- البروتينات ذات الصلة.
- تحورات ما بعد الترجمة.
- التفاعلات مع البروتينات الأخرى.
- روابط اتصال بقواعد بيانات أخرى.

وقد ساهمت البيومعلوماتية فى انشاء وتطوير قواعد البيانات هذه، وكذلك ايجاد نظم خوارزمية لمقارنة وتحليل طرز البروتينات المحتوية عليها تلك القواعد.

قواعد بيانات المسارات الأيضية **Databases of metabolic pathways**:

تقوم قاعدة KEGG (Kyoto Encyclopedia of Gens and Genomes) بجمع الجينوم، ونواتج الجين ووظائفها، وعمل تكامل بين المعلومات البيوكيميائية والوراثية. وترتكز KEGG على تفاعلات شبكات التجميع الجزيئى والشبكات الأيضية والمنظمة.

(١) وتنظم KEGG خمسة أنواع من البيانات من خلال نظام شامل وهى:

- (٢) كتالوج للجينات يحتوى على جزيئات وتتابعات معينة.
- (٣) خرائط الجينوم. حيث يتم التكامل بين الجينات طبقا لظهورهم على الكروموسومات.
- (٤) خرائط المسارات والتي تصف شبكات الأنشطة الأيضية والمنظمة للجزيئات. وتساهم فى امكانية ايجاد مسار ايسى حقيقى فى كائن معين عن طريق مضاهاة بروتينات هذا الكائن مع انزيمات مسارات مرجعية.
- (٥) جداول أرثولوج. تستخدم بيانات تلك الجداول لربط انزيم فى كائن ما بالانزيمات ذات الصلة فى كائنات أخرى. ويتيح ذلك تحليل العلاقات بين المسارات الأيضية فى كائنات مختلفة.
- تتيح KEGG امكانية أخذ مجموعة من الانزيمات من كائن واختبار مدى تكاملها مع مسارات ايسية معروفة. وفى حالة ظهور فراغ فى مسار ما فان ذلك قد يرجع الى غياب انزيم أو مسار ايسى غير متوقع.

مداخل الى السجلات :

تمتلك قواعد بيانات الحمض النووى وتتابعات البروتين امكانات هائلة لاسترجاع المعلومات وتحليلها من خلال عديد من العمليات والتي تتضمن:

١. استرجاع التتابعات من قاعدة البيانات.
٢. مقارنة التتابع.
٣. ترجمة تتابعات الدنا الى تتابعات بروتين.
٤. تحليل التركيب والتنبؤ.
٥. طرز التعرف.
٦. الأشكال الجزيئية. وتستخدم تلك الأشكال فى:

- خرطنة الأجزاء التي يعتقد أن لها وظيفة ما الى اطار ثلاثى الأبعاد لبروتين.
- تقسيم ومقارنة طرز الطى فى البروتينات.
- تحليل الاختلافات بين التراكيب شديدة القرابة أو بين الأوضاع الفراغية لجزئ ما.
- دراسة تفاعل جزئ بسيط مع بروتين كمحاولة لتحديد وظيفة أو لاكتشاف دواء.
- الحصول على نموذج تفاعلى لتحسين جودة صور الأشعة السينية لمقاييس تراكيب البروتين.
- تصميم ونمذجة تراكيب جديدة.

ENTREZ: هى نقطة البداية لاسترجاع وتتابعات وتراكيب من على الموقع التالى:

<http://www.ncbi.nlm.nih.gov/ENTREZ/>

طرق الدخول إلى قواعد
بيانات البيولوجيا
الجزئية:

وتقدم ENTREZ طرق دخول عبر الأقسام الآتية لقواعد البيانات:

- Protein
- Peptide
- Nucleotide
- Structure
- Genome
- Popset- information about populations
- OMIM- Online Mendelian Inheritance in Man

فمثلا للبحث في قاعدة البروتين: ادخل على موقع ENTREZ، ثم اختر بروتين، ثم ادخل نوع البروتين المطلوب البحث عنه، ثم اضغط على GO وسيقوم البرنامج باظهار الاجابات.

وللبحث في قاعدة بيانات النيكلوتيدات: اختر نيكلوتيد من على موقع ENTREZ واضغط على INDEX ثم اختر الكائن من القائمة وادخل اسم الموضوع في صندوق البحث، ثم اضغط على AND ثم اتبع التعليمات من القائمة.

:OMIM

وهي قاعدة بيانات للجينات البشرية والعيوب الوراثية. وتعمل بتكامل تام مع ENTREZ.

نظام استرجاع التتابع (SRS) The Sequence Retrieval System:

يستطيع نظام SRS البحث في ١٤١ من قواعد بيانات البروتين، وتتابعات النيكلوتيد، والمسارات الأيضية، والتراكيب ثلاثية الأبعاد، والوظائف، والجينوم، والأمراض.

مصدر تعريف البروتين (PIR) The Protein Identification Resource:

يمكن البحث باستخدام SRS من على المواقع التالية :
في الولايات المتحدة الأمريكية :

<http://www.-nbrf.georgetown.edu/pirwww/search/textpsd.html>

في أوروبا: <http://www.mips.gsf.de>

: ExPASy-Expert Protein Analysis System

وهو نظام استرجاع وتحليل المعلومات لمعهد البيومعلوماتية السويسري بالتعاون مع معهد البيومعلوماتية الأوروبي والذي ينتج قواعد بيانات تتابع البروتين SWISS-PROT وقاعدة TrEMBL التي تحتوى على تراجم تتابعات النيكلوتيد. وفتح الصفحة الرئيسية ل ExPASy

(<http://www.expasy.ch>) واختيار SWISS-PROT و TrEMBL يعطى
مداخل توصل الى أدوات استرجاع المعلومات.

:ENSEMBL

وهو عبارة عن مصدر معلومات شامل للجينوم البشري
(<http://www.ensembl.org>). وتهدف الى جمع وتفسير كل المعلومات
المتاحة عن تتابعات الدنا البشري، وروابط اتصال التي تتابع جينوم
قياسي، وإتاحة هذه المعلومات لعدد من العلماء. ويمكن للمستخدم تعريف
مناطق في التتابع عبر عدة أدوات بحث:

- BLAST ويستخدم للبحث عن تتابع أو جزء منه.
- التصفح، ويبدأ على مستوى الكروموسوم ثم يتعمق داخله.
- العلاقة بالأمراض بواسطة OMIM.
- ENSEMBL ID.
- أداة بحث عامة للنصوص.

الفصل الرابع

المحاذاة والاشجار العرقية

محاذاة التتابع:

في حالة وجود اثنان من التتابعات يتم عمل الآتي:

- قياس درجة التشابه بينهما.
- تحديد تقابل الأجزاء.
- ملاحظة طرز التحفظ والاختلاف.
- علاقات النشوء.

:The dotplot

وهي صورة بسيطة تعطي رؤية عامة للتشابه بين اثنان من التتابعات. يتكون dotplot من جدول أو مصفوفة. تشير الصفوف الى أجزاء أحد التتابعات بينما تدل الأعمدة على التتابع الآخر. وتظهر بقع في حالة التماثل، بينما لا توجد استجابة في حالة الاختلاف.

مصادر DOTPLOTS على الانترنت (web Resources):

يتيح برنامج E.L. Sonnhammer's Dotter للمستخدم التحكم وطرق الحساب وتغيير مظهر العرض بواسطة ضبط معايير تفاعلية والذي يجب تجهيزه في الكمبيوتر من على الموقع التالي:

<http://www.cgr.ki.se/cgr/groups/sonnhammer/Dotter.html>

يقدم الموقع التالي امكانية اجراء dotplotting تفاعلي:

<http://www.isrec.isb-sib.ch/java/dotlet/exonintron.html>

A PERL program to draw dotplots.

```
#!/usr/bin/perl
#dotplot.pl -- reads two sequences and prints dotplot
# read input
$/ = "";
$_ = <DATA> $_ =~ s/#!/.\n\n/g;
$_ =
    =~
    /^(.*)\n\s*(\d+)\s+(\d+)\s*\n(.*)\n([A-
Z\n]*)\s*\n(.*)\n([A-Z\n]*)\s*/;
$title = $1; $nwind = $2; $thresh = $3;
$seqt1 = $4; $seq1 = $5; $seqt2 = $6; $seq2 = $7;
$seq1 =~ s/\n//g; $seq2 =~ s/\n//g; $n = length($seq1); $m =
length($seq2);
# postscript header
print << EOF setfont scalefont 20 findfont Helvetica translate 30
setlinewidth 1.75 def load fill f closepath c newpath n rlineto r
moveto m lineto l stroke s !PS-Adobe->

#print matrix

$dx = 500.0/$n; $mdx = -$dx; $dy = 500.0/$m;
if ($dy < $dx) {$dx = $dy;} $dy = $dx; $xmx = $n*$dx; $ymx =
$m*$dx;
print "0 510 m ($title NWIND = $nwind) show\n";
printf "0 0 m 0 %9.2f 1 %9.2f %9.2f 1 %9.2f 0 1 c s\n",
$ymx, $xmx, $ymx, $xmx;
for ($k = $nwind - $m + 1; $k < $n - $nwind; $k++) {
    $i = $k; $j = 1; if ($k < 1) {$i = 1; $j = 2 - $k;}
    while ($i <= $n - $nwind && $j <= $m - $nwind) {
        $_ = (substr($seq1, $i - 1, $nwind) ^ substr($seq2, $j -
1, $nwind));
        $mismatch = ($_ =~ s/^\x0//g);
        if ($mismatch < $thresh) {
            $xl = ($i - 1)*$dx; $yb = ($m - $j)*$dy;
            printf "n %9.2f %9.2f m %9.2f 0 r 0 %9.2f r %9.2f 0 r c
f\n",
                $xl, $yb, $dx, $dy, $mdx;
        }
        $i++; $j++;
    }
}
print "showpage\n";

__END__
ATPases lamprey / dogfish #TITLE
15 6 #WINDOW, TRESHOLD
Petromyzon marinus mitochondrion #SEQUENC : 1
ATGACACTAGATATCTTTGACCAATTTACCTCCCAACA
ATATTTGGGCTTCCACTAGCCTGATTAGCTATACTAGCCCCTAGCTTA
ATATTAGTTTCACAAACACCAAATTTATCAAATCTCGTTATCACACACTA
CTTACACCCATCTTAACATCTATGCAAAACAACCTCTTTCTTCCAATAAAC
CAACAAGGCATAAATGAGCCTTAATTTGTATAGCCTCATAATATTTATC
TTAATAATTAATCTTTTAGGATTATTACCATATACTTATACACCAACTACC
CAATTATCAATAAACATAGGATTAGCAGTGCCACTATGACTAGCTACTGTC
CTCATTGGGTTACAAAAAACAACAGAGCCCTAGCCCCTATTATACCA
GAAGGTACCCAGCAGCACTCATCCCATATTAATTATCATTTGAAACTATT
AGTCTTTTATCCGACCTATCGCCCTAGGAGTCCGACTAACCGCTAATTTA
ACAGCTGGTCACTTACTTATACAACCTAGTTCTATAACAACCTTTGTAATA
ATTCTGTCAATTTCAATTTCAATTTACTCTACTACTCTCTCTATTA
CTAACAATTTCTGGAGTTAGCTGTTGCTGTAATCCAGGCATATGATTTATT
CTACTTTTAACTCTTTATCTGCAAGAAAACGTTT*
Scyliorhinus canicula mitochondrion #SEQUENCE 2
ATGATTATAAGCTTTTTGATCAATTCCTAAGTCCCTCCTTTCTAGGA
ATCCCACTAATGGCCTAGCTATTTCAATTCATGATTAATTTCCAACACCAAC
AATCGTTGACTTAAATAATCGATTATTAACCTTCAAGCATGATTTATTAACCGATT
TATCAACTAATACAACCCATAAATTTAGGAGGACATAAATGAGCTATCTATTTACA
CTAATATTTATTTAATTTACCATCAATCTTCTAGGTCTCCTTCCATATACTTTTACG
CT
ACAACTCAACTTTCTTAAATATAGCCTTTGCCCTGCCCTTATGGCTTACAACCTGATTA
```

Copyright © 2006. All rights reserved. May not be reproduced in any form without permission from the publisher, except fair uses permitted under U.S. or applicable copyright law.

```
ATTGGTATATTTAATCAACCAACCATTTGCCCTAGGGCACTTATTACCTGAAGGTACCCCA  
ACCCCTTTAGTACCAGTACTAATCATTATCGAAACCATCAGTTATTTATTCGACCATTA  
GCTTAGGAGTCCGATTAACAGCCAACCTAACAGCTGGACATCTCCTTATACAATTAATC  
GCAACTGCGGCCTTTGTCCTTTTAACTATAAATACCAACCGTGGCCTTACTAACCTCCCTA  
GTCCTGTTCCATATGACTATTTTAGAAGTGGCTGTAGCTATAATTCAAGCATACTATTT  
GTCCTTCTTTTAAAGCTTATATCTACAAGAAAACGTATAA*
```

محازاة التتابع بواسطة Dotplots:

لا يقدم Dotplots فقط صورة شاملة لتشابه اثنان من التتابعات ولكن أيضا يعطى مجموعة متكاملة ذات جودة نسبية للمحازاة المحتملة وذلك بالمرور على الصورة من الجانب الأيسر العلوى الى الجانب الأيمن السفلى.

محازاة التتابع المتعدد Multiple sequence alignment :

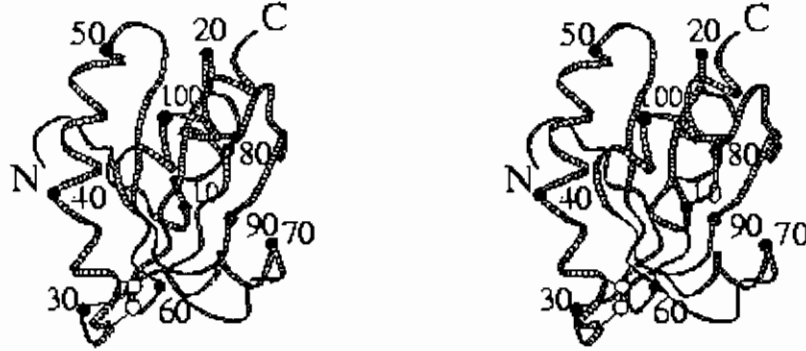
تكون نتائج الكشف عن العلاقات البعيدة بين التتابعات أكثر واقعية فى حالة تعدد التتابعات. وأيضا تعطى أدوات التنبؤ بالتركيب نتائج أفضل عندما تعتمد على محازاة تتابعات متعددة وذلك بالمقارنة بالتتابعات المنفردة. ولذلك يجب أن تحتوى محازاة التتابع المتعدد على توزيع من التتابعات قريبة وبعيدة الصلة.

الاستدلال التركيبى بواسطة محازاة التتابع المتعدد:

يمكن التعرف على الملامح التركيبية والوظيفية بواسطة المحازاة للتتابعات المتعددة من خلال موقع الخادم التالى:

<http://www.bio.cam.ac.uk/cgi-bin/seqlogo/logo.cgi>

ومثال ذلك التعرف على الملامح التركيبية لعدد ١٦ ثيوريدوكسينات (انزيمات موجودة فى كل الخلايا وتساهم فى مدى واسع من العمليات البيولوجية):



The structure of *E. coli* thioredoxin [2TRX] contains a central five-stranded β -sheet flanked on either side by α -helices. Residue numbers correspond to those in the multiple sequence alignment table. The N- and C-termini are also marked. Spheres indicate positions of the $C\alpha$ atoms of every tenth residue. The reactive disulphide bridge between Cys32 and Cys35 appears between the numerals 30 and 60 .

Hidden Markov Models (HMMS) :

وهي نماذج تركيب حسابية لوصف طرز عائلات التتابعات المتماثلة. كما أنها تعتبر من الأدوات القوية للكشف عن القرابة البعيدة والتنبؤ بطرز طي البروتين.

مصادر على الانترنت : Web Resources

<http://cse.ucsc.edu/research/compbio/sam.html>

<http://cse.ucsc.edu/research/compbio/HMM-apps/HMM-applications.html>

<http://www.sanger.ac.uk/Software/Pfam>.

<http://stash.mre-lmb.cam.ac.uk/SUPERFAMILY/>

النشوء النوعي : Phylogeny

هناك العديد من الأمثلة للتطور في البروتينات والجينوم. ويهدف هذا المجال العمل على التوصل الى العلاقات بين الأنواع والعشائر والأفراد والجينات وذلك باستخدام الخصائص الجزيئية.

وتمثل الشجرة في علم الحاسبات شكلا خاصا يتكون من نقط nodes يتم توصيلها ببعض عن طريق خطوط edges. والشجرة قد تكون rooted أو unrooted. والشكل الأخير لا يظهر النسب بينما في الشكل الأول ينحدر من كل نقطة أصلان ويسمى binary tree ويستخدم برنامج PERL لرسم الشجرة كما يلي:

drawtree.prl - PERL

program to draw binary trees.

```
#!/usr/bin/perl
#drawtree.pl -- draws binary trees (root at top)
#usage: echo '(A((BC)D)(EF))' | drawtree.pl > output.ps

print <; chop($tree); $_ = reverse($tree); s/[()]/g;

$x = 0; $y = 0;
while ($nd = chop()) {
    print "$x $y m ($nd) stringwidth pop -0.5 mul 0 rm ($nd)
show\n";
    $xx{$nd} = $x; $x+=20; $yy{$nd} = 10;
}

while ($tree =~ s/(?{[A-Z]}{[A-Z]}?/1/) {
    print "n $xx{$1} $yy{$1} m\n";
    ($yy{$1} > $yy{$2}) || ($yy{$1} = $yy{$2}); $yy{$1} += 20;
    print "$xx{$1} $yy{$1} l $xx{$2} $yy{$1} l $xx{$2} $yy{$2} l
s\n";
    $xx{$1} = 0.5*($xx{$1} + $xx{$2});
}print "n $xx{$tree} $yy{$tree} m 0 20 rl s showpage\n";

$rx = 2*$x + 30; $yt = 2*$yy{$tree} + 146;
print "%BoundingBox: 40 95 $rx $yt\n";
```

مصادر على الانترنت لأدوات الوراثة العرقية :

<http://evolution.genetics.washington.edu/phylip/software.html>

الفصل الخامس

تركيب البروتين واكتشاف الدواء

يتضمن هذا الفصل من الكتاب تركيب البروتينات من حيث:

يستخدم تحليل hydrophobicity profile للتنبؤ بأماكن التفاف العناصر في التركيب الثانوي، الوحدات الظاهرة والمخفية، امتداد أجزاء الغشاء، ومواقع الأنتيجين. وفيمايلي مثال لاستخدام برنامج PERL لرسم مايسمى بـ The helical wheel التي تعتمد على دراسة تتابع الوحدات من حيث التبادل ما بين محب وكاره للماء وذلك في alf-helics in globular proteins

**ثبات وطى البروتين
وتطبيقه**
:Hydrophobicity

helwheel.pl - program to draw helical wheels.

```
#!/usr/bin/perl
#helwheel.pl -- draw helical wheel
#usage: echo DVAGHGQDILIRLFKSH | helwheel.pl > output.ps
# or   echo 20DVAGHGQDILIRLFKSH | helwheel.pl > output.ps
#       the numerical prefix sets the first residue number

# The output of this program is in PostScript (TM),
#       a general-purpose graphical language

# The next section prints a header for the PostScript file

print <;                                # read line of input
chop();$_ = s/\s//g;                     # remove terminal
carriage return and blanks

if ($_ =~ s/^(\\d+)/)                    # if input begins with
integer                                  #
{ $resno = $1;                            # extract it as
initial residue number
} else { $resno = 1;                       # if not, set initial
residue number = 1

$radius = 50;                             # initialize values
for radius,
{x = 0; $y = -50; $theta = -90;          # x, y and angle theta
```

```

# print light gray spiral arc as succession of line segments, 10
per residue

$npoints = 10*(length($_) - 1);
print "0.8 0.8 0.8 setrgbcolor\n"; # set colour to light
gray
print "newpath\n"; # draw spiral arc
printf("%8.3f %8.3f moveto\n", $x, $y);
foreach $d (1 .. $npoints) { # 10 points per
residue
    $theta += 10; $radius += 0.6; # increase radius and
theta
    $x = $radius*cos($theta*0.01747737); # calculate new value
of x
    $y = $radius*sin($theta*0.01747737); # and y
    printf("%8.3f %8.3f lineto\n", $x, $y);
}
print "stroke\n";

# print residues and residue numbers

$radius = 50; # reinitialize values
for radius,
$x = 0; $y = -50; $theta = -90; # x, y and angle theta
print ``0 setgray\n`` # set colour to black

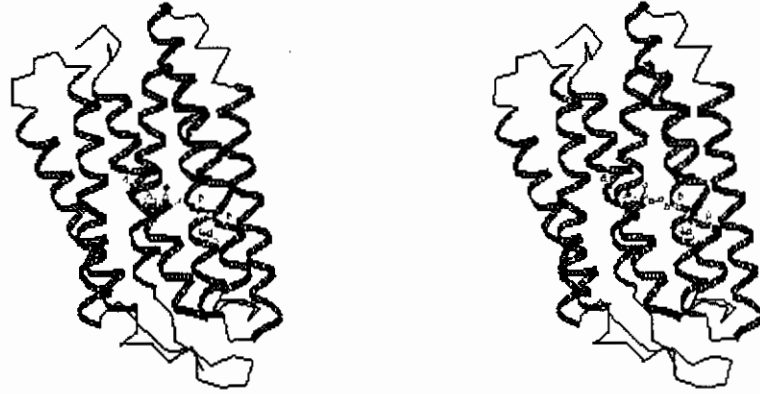
foreach (split ("", $_)) { # loop over characters
from input line
    print "/$font($_) findfont"; # set font appropriate
    print "20 scalefont setfont\n"; # for this amino acid
    printf("%8.3f %8.3f moveto\n", $x, $y); # move to current
point
    print " ($resno$_) stringwidth"; # adjust position to
center residue
    print " pop -0.5 mul -7 rmoveto\n"; # identification on
point on spiral
    print " ($resno$_) show\n"; # print residue number
and id
    print "% $theta $resno$_\n";
    $theta += 100; $radius += 6; # set new values of
angle, radius
    $x = $radius*cos($theta*0.01747737); # compute new values
of x
    $y = $radius*sin($theta*0.01747737); # and y
    $resno++; # increase residue
number
}

print "showpage\n"; # postscript signals
to
print "%BoundingBox:"; # print
$x1 = 297.5 - 1.05*$radius; # x
$x2 = 297.5 + 1.05*$radius; # and
$y1 = 421. - 1.05*$radius; # y
$y2 = 421. + 1.05*$radius; # limits
printf("%8.3f %8.3f %8.3f %8.3f\n", $x1, $x2, $y1, $y2);

print "showpage\n";
print "%EOF\n"; # and wind up

```

المثال الثاني: التنبؤ بالتركيب اللولبي لبروتينات الأغشية، حيث أن العديد من بروتينات الأغشية تتركب من سبع أشكال لولبية متصلة مع بعضها بواسطة عقد كما هو الحال في bacteriorhodopsin :



Bacteriorhodopsin membrane from the bacterium *Halobacterium salinarum* . The ligand shown in ball-and-stick representation is the chromophore, retinal .

Web Resources : Transmembrane Helix Prediction :

TMHMM (A.Krogh and E. Sonnhammer
Markov Model :
<http://www.cbs.dtu.dk/krogh/TMHMM/>

PHDhtm (B.Rost) :
<http://dodo.bioc.columbia.edu/predictprotein>

Membrane protein explorer (S. White) :
<http://blanco.biomol.uci.edu/mpex/>

Web Resources :

matrix Alignment) :
<http://www2.ebi.ac.uk/dali/>

Classifications of Protein Structure :
<http://www2.ebi.ac.uk/msd/Links/fold.shtml>.

<http://www.bioscience.org/urlists/protdb.htm>
<http://scop.mrc-lmb.cam.ac.uk/scop>

Protein Structure Prediction and modelling :
<http://cubic.bioc.columbia.edu/eva>

Homology Modelling :

SWISS-MODEL (automatic homology modelling) :
<http://www.expasy.ch/swissmod/SWISS-MODEL.html>

MODBASE, a database of comparative models of protein from complete genomes:
<http://guitar.rockefeller.edu/modbase/>

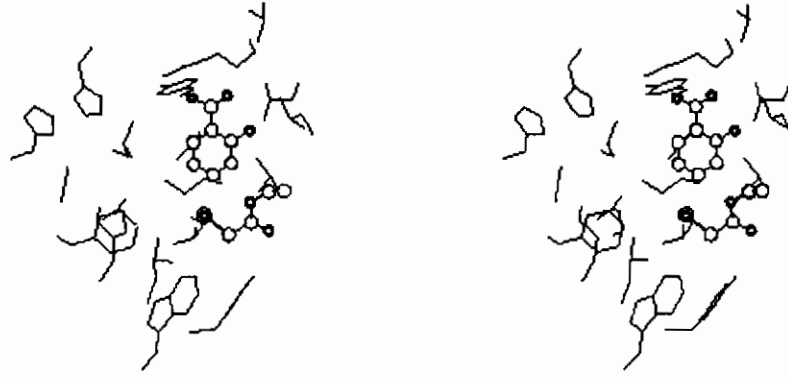
تم استعراض المواصفات التي يجب توافرها في المركب الكيميائي المؤهلة لاستخدامه كدواء. وكذلك مراحل انتاج دواء جديد. واستخدام الكمبيوتر في حصر مركبات عديدة للوصول الى مركب واحد مؤهل، ثم اجراء تحسين للمركب المرشح من خلال تحديد العلاقة الكمية بين التركيب والفعالية QSAR والتي تتيح طرق للتنبؤ بالنشاط الدوائي لمجموعة من المركبات بواسطة العلاقة بين الملامح الجزيئية والنشاط الدوائي. وكذلك الاستعانة بمعلومات عن الجينوم.

اكتشاف وتطوير الدواء:

تصميم الدواء بمساعدة الكمبيوتر:

مثال: تصميم مثبطات 2 prostaglandin cyclo-oxygenase. ومن المعروف أن مركب الأسبرين ومركبات أخرى من NSAID تثبط اثنان من prostaglandin cyclo-oxygenase قريبي الصلة تسمى COX-1 و COX-2

: 2



The binding site in COX-1 for an aspirin analog, 2-bromoacetoxybenzoic acid . The ligand has reacted with the protein, transferring the bromoacetyl group to the sidechain of serine 530 . The protein is shown in skeletalrepresentation . The aspirin analog is shown in ball-and-stick representation .

وللمزيد من المعلومات يمكن الرجوع الى الموقع الالكتروني للكتاب:

<http://www.oup.com/uk/lesk/bioinf>

الخاتمة

استقراء الحالة الراهنة لتطبيقات البيومعلوماتية تمكننا من استشراف أهمية مجال البيومعلوماتية في المستقبل. فمن الواضح أن عمليات جمع البيانات ستستمر في الزيادة وبسرعة. كذلك التوجه نحو زيادة كفاءة وقدرة أجهزة الكمبيوتر من حيث التخزين، والتوزيع، وتحليل النتائج. وسوف يؤدي ابتكار طرق خوارزمية محسنة الى تحليل وتفسير المعلومات المتاحة لنا وتحويلها من مجرد صور بيانات الى معرفة ثم حكمة.

وسوف نصل الى الكتلة الحرجة عندما تصبح معرفتنا بالتتابعات والتراكيب أكثر قربا من الاكتمال، بمعنى جمع مجموعات غير منحازة من البيانات المتاحة عن أشكال الحياة المعاصرة. وسوف نتعرف على ذلك عندما يؤدي التعمق العشوائي في جينوم ما أو معرفة تركيب بروتين ما الى تغيير شيئا موجودا بالفعل بدلا من الكشف عن شيئا جديدا. وبعد كل ذلك، تتكون الطبيعة من نظام غير محدود الاحتمالات ولاكن مع اختيارات لانهائية.

ستصبح التطبيقات أكثر وضوحا، ونضجا وستمر بسرعة من مراحل البحث الى الممارسات الصناعية والاكاديمية القياسية. ونقل بعض المعلومات البيولوجية عالية المستوى مثل برامج التطور أثناء فترة حياة الأفراد، وأنشطة العقل البشري ودخولها في عمليات معالجة البيانات، سوف يمكننا من الوصف الكمي والتحليل على مستوى الجزيئات وتفاعلها.

رقم الإيداع: ٢٠٠٦/٢١٩١

ISBN: 977-281-297-5

مطابع الجدار الهندسية/القاهرة
تليفون/فاكس : (٢٠٢) ٥٤٠٢٥٩٨