

کارآیی نامعقول ریاضیات در زیست‌شناسی مولکولی*

آرتور لسک*

ترجمه افرا علیشاھی

شیمی نتایج ساده‌ای به دنبال داشته باشد که توصیفات مبسوطی از فرآیندهای حیات در اختیار ما قرار دهد، ما ممکن است قادر به کشف آنها نباشیم، زیرا پدیده‌های مورد مطالعه ما پیچیده‌ترند، در برابر ایده‌آل‌سازی‌های ساده‌گذرنده مقاومت می‌کنند، و خصوصیاتی از خود نشان می‌دهند که تحت تأثیر شدید انتخاب شرایط اولیه از بین مجموعه بسیار عظیم و پراکنده‌ای از شرایط ممکن است. سیهای در زیست‌شناسی نقشی بسیار مهمتر از افتادن بر سر مردم ایفا می‌کنند.

موضوع بحث در زیست‌شناسی مولکولی

ایشای مورد مطالعه ما حداقل شکلی دارند که می‌توانیم سعی کثیم ریاضیات را در مورد آنها به کار گیریم. این اشیاء عبارت‌اند از:

- توالیهای زناه در DNA
- توالیهای آمینواسیدها در پروتئینها
- ساختار پروتئینها
- کارکرد پروتئینها

احتمالاً خوشنودگان در باره پروژه «نوم شنیده‌اند» هدف این پروژه تعیین توالی کامل DNA در موجودات زنده است: یعنی مجموعه نقشه‌ها. توالیهای DNA در زونهای حاوی همه اطلاعاتی هستند که یک موجود زنده برای تولد، بزرگ شدن و رشد، و مرگ به آنها نیاز دارد. با تکمیل تعیین توالی زنوم مخمر در سال ۱۹۹۶، ما همانقدر در باره یک سلول مخمر می‌دانیم که خود سلول مخمر می‌داند. این عبارت آنقدر که به نظر می‌رسد متکبرانه نیست. ما واقعاً همه اطلاعات را داریم. باید پذیرفت که ما نمی‌توانیم این اطلاعات را به همان کارآمدی تفسیر کنیم که یک سلول مخمر می‌تواند، اما ما مجموعه کامل نقشه‌ها را در اختیار داریم. ولی این نقشه‌ها تنها یک توصیف است از ساختار و عملکرد بالقوه بدست می‌دهند؛ می‌ماند اینکه مشاهدات خود را تا مرحله یکپارچه‌سازی توصیف و کارکرد پروتئین در زمان و فضای درون یک موجود

عنوان مقاله من اقتباسی است از عنوان مقاله مشهور ویگنر، «کارآیی نامعقول ریاضیات در علوم طبیعی [۱]». البته این عنوان در فیزیک و در زیست‌شناسی مولکولی، از دو بابت مخالف توجه برانگیز است. در فیزیک، بدیهی است که ریاضیات کارآمد است — بسیاری از غولهایی که فیزیکدانها بر شانه‌های آنها ایستاده‌اند ریاضیداناند — و مایه شگفتی است که ویگنر را نامعقول می‌داند. در زیست‌شناسی مولکولی، نقش واقعی ریاضیات بدیهی نیست، و بیم آن می‌رود که انتظار کارآیی از ریاضیات نامعقول باشد، بیمی که در این مورد از فیزیک بسیار موجه‌تر است. البته، بسیاری از ایده‌های متداوی در زیست‌شناسی مولکولی محاسباتی — مثلاً جستجو در پایگاه داده‌ها برای یافتن دنباله‌های شبیه به یک دنباله یافت شده — مسلماً مبتنی بر ریاضیات و علوم کامپیوتر هستند. اما اینکه آیا درک غایی ما از فرآیندهای حیاتی به زبان ریاضیات بیان شود — آن چنان‌که، مثلاً مفاهیم تقارن زمینه‌ساز بیان قوانین فیزیک هستند — یا به زبان سنتی توصیفی «روایی» زیست‌شناسی، هنوز محل بحث است.

چرا تردید در کارآیی ریاضیات در زیست‌شناسی می‌تواند معقول باشد؟ خصوصیات مشاهده‌شده سیستمهای زنده به وسیله ترکیبی از اینها تعیین می‌شوند:

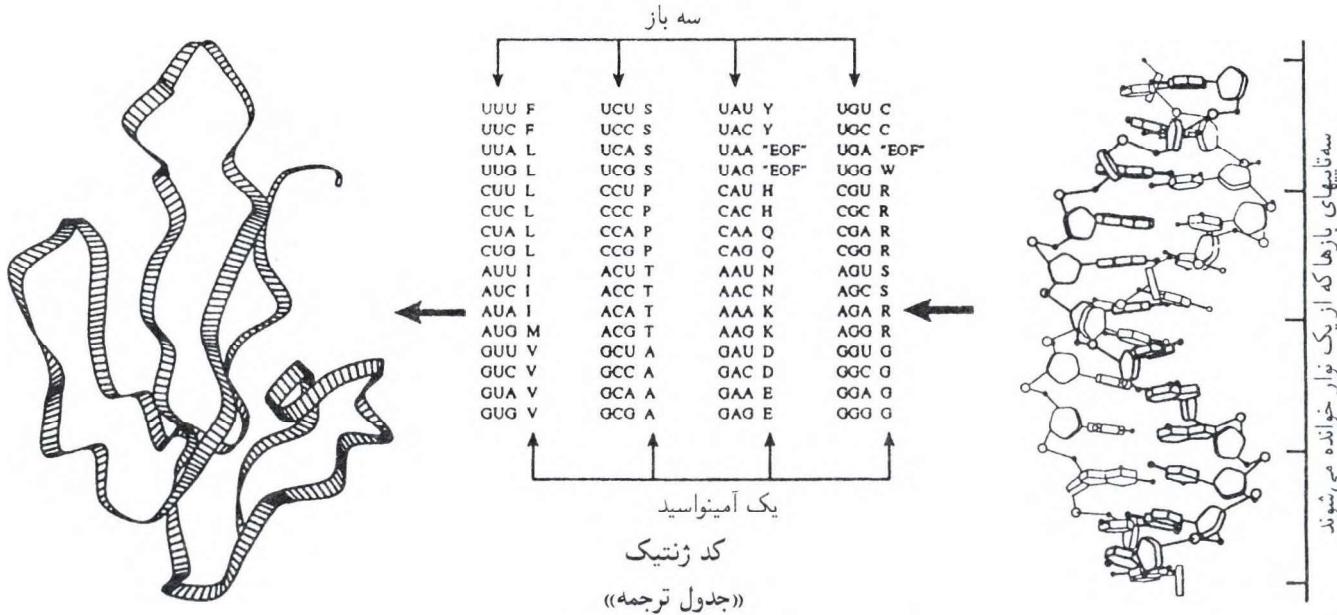
- قوانین فیزیک و شیمی
- سازوکار تکامل
- تصادف تاریخی

تفکیک تأثیر اینها از هم دشوار است، و کشاکش سازنده‌ای که بین آنهاست، تحقیقات ما را تحت تأثیر قرار می‌دهد. بسیاری از قوانین فیزیک، جهان طبیعی را — که شامل سیستمهای زنده هم می‌شود — با تعیین روابط بین شرایط اولیه و نهایی توصیف می‌کنند. در زیست‌شناسی پیچیدگی مجموعه شرایط اولیه ممکن، دشواری‌هایی پدید می‌آورد. نقش عظیم تصادف تاریخی ما را به تردید می‌اندازد و به زانو درمی‌آورد: حتی اگر قوانین بنیادی فیزیک و

که خود به خود تا می‌شود تا
ساختار سه‌بعدی دقیقی بسازد

به توالی‌ای از آمینواسیدها در
پروتئین ترجمه می‌شود ...

... توالی‌ای از بازها در DNA



شکل ۱ جریان اطلاعات در خلال «خوانش» از روی یک ژن. ژنهای، یا نقشه‌های بنیادی موجودات زنده، در ساختار DNA گنجانده شده‌اند (راست). در سمت راست شکل، مارپیچ مضاعف DNA را می‌بینیم، که حاوی دو نوار بهم پیچیده است که یکی با خطوط باریک و دیگری با خطوط پهن رسم شده است. (اتر پلکانی) توسط مجموعه‌ای از زیرواحدهای شیمیایی به نام «باز» پدید می‌آید. در هر نقطه روی هر نوار یکی از چهار باز A, T, C, G یا C ممکن است قرار گیرند. خواننده تیزین می‌تواند بینند که این بازها — «لبه‌های قائم» پله‌های پلکان مارپیچی — به شکلهای مختلف ظاهر می‌شوند. بازهای همسطح هر دو نوار با هم ارتباط برقرار می‌کنند، و این ارتباط مستلزم مکمل بودن طبق است: حضور یک A روی یک نوار مستلزم حضور یک T روی نوار دیگر است، و یک G روی یک نوار به یک C روی نوار دیگر نیاز دارد، و برعکس، یک T روی یک نوار مکمل یک A و یک C مکمل یک G روی نوار دیگر است. به این ترتیب، هر نوار حاوی اطلاعات کافی برای ساختن نوار همراه خود است. به لحاظ منطقی، راه تکثیر DNA آن است که نوارها را از هم جدا کنیم و مکمل هر کدام از نوارهای مجزا را دوباره پسازیم. توالیهای جازها در ژنهای، همکمل یک جدول ترجمه مستقیم عذام «کد ژنتیک» توالیهای آمینواسید دو پروتئینها (هرگشایی می‌کنند (وسط)). پیام ژنتیکی با الفای چهار حرفی A, T, C, G و نوشته می‌شود. پروتئینها هم پایم‌هایی هستند که توالی‌ای از سازه‌های شیمیایی را در خود دارند، چنان‌که در مر نطفه یکی از بیست آمینواسید ممکن قرار می‌گیرد، که با نمادهای A, T, C, G, F, E, D, C, A, W, V, T, S, R, Q, P, N, M, L, K, H, G, F, E, D, C, A مشخص می‌شوند. برای مشخص کردن مجموعه ۲۰ آمینواسید، هر کدام به بیش از دو باز نیاز دارد؛ درواقع از روی توالی DNA در یک زمان سه باز خوانده می‌شود، با افزونگی‌ای که برای تکامل مطلقاً ضروری است (سه سمت‌ای از بازها برای علامتهای انتهایی «یابان حیات» نگهداری شده‌اند). پروتئینها خود به خود تا می‌شوند تا ساختارهای طبیعی فعال سه‌بعدی بسازند (چپ). این تقطه‌ای است که در آن طبیعت از توالیهای یک بعدی ژنهای به جهان سه‌بعدی ای که ما در آن زندگی می‌کنند. مثال فوق یک توکسین از یک مار آبی است، یکی از ساختارهای پروتئین متعددی که به وسیله بلورنگاری با اشعه X تعیین شده است. هر ژن توالی‌ای از بازها دارد، که ابتدا توالی آمینواسیدهای یک پروتئین، و سپس ساختار سه‌بعدی آن، و سپس کارکرد آن را تعیین می‌کند.

(شکل ۱). برای تعیین کارکرد پروتئین، یک ساختار سه‌بعدی دقیق ضروری است، زیرا تعاملهای مورد نیاز به کنار هم قرار گرفتن بخش‌های مختلف مولکولها در روابط فضایی دقیق بستگی دارند. نهایتاً پسخورد کارکرد پروتئین به دنباله ژنهای — از طریق تکامل به روش انتخاب طبیعی — این حلقه را کامل می‌کند. توالیهای DNA در کامپوترهای ما، رشته‌های حرفي هستند، یعنی اشیاء یک بعدی. ژنهای یا زیررشته‌هایی از توالیهای ژنوم، براساس یک کدنویسی تقریباً جهان‌شمول به توالیهای آمینواسید پروتئینها ترجمه می‌شوند. این توالیها هم به وسیله رشته‌های حرفي یک بعدی بازنمایی می‌شوند. سپس پروتئینها خود به خود تا می‌شوند تا ساختارهای سه‌بعدی «طبیعی» یکتا سازند. (شاهد این مدعای آن است که می‌توان آنها را به وسیله حرارت تحریب کرد — ساختار سه‌بعدی را ویران کرد — و وقتی سرد شوند مجدداً شکل اولیه خود را پیدا می‌کنند، مانند آلیازهای شکل-حافظه^۱). تاخوردگی خود به خود پروتئینها

زنده بسط دهیم. جمع‌آوری این اطلاعات به «پروژه پروتوم» معروف است. پروژه‌ای که در حال جمع‌آوری قواست تا در عصر پس از زنوم ایفای نقش کند. میزان اندازه‌گیری توالیهای ژنهای بسیار زیاد و فزینده است. در سال ۱۹۹۸ توالی کامل DNA^۱ یک نوع کرم موسوم به Caenorhabditis elegans تکمیل شد (۹۷×۱۰^۶ باز). محتمل است که در سالهای ۱۹۹۹ و ۲۰۰۰ به ترتیب شاهد تکمیل توالی DNA^۲ می‌گرسی سرمه (۱×۱۰^۴ باز) و ژنوم انسان (۳×۱۰^۹ باز) و بسیاری موجودات زنده کوچک و بزرگ دیگر باشیم. لویی پانزدهم می‌توانست بگوید: «بعد از من، طوفان»، اما نوح نمی‌توانست: ما هم نمی‌توانیم.

توالیها و ساختارهایی که ما مطالعه می‌کنیم ارتباطات مهمی با هم دارند. در سطح مولکول، توالیهای ژنهای DNA، توالیهای آمینواسیدهای پروتئینها را رمزگشایی می‌کنند. سپس توالیهای آمینواسید پروتئینها ساختار سه‌بعدی پروتئینها را تعیین می‌کنند.

چگونه می‌توان تعیین کرد که کدام‌یک از این دو، یا کدام‌یک از بسیاری از هم‌ردیف‌سازی‌های ممکن دیگر، بهترین هم‌ردیف‌سازی است؟ آیا ما می‌توانیم متريکی برای رشته‌های حرفی طرح کنیم و فاصله بین آنها را تعریف کنیم؟

شاخصهای عدم شباهت بین رشته‌های حرفی عبارت‌اند از:

(۱) فاصله همینگ، که بین دو رشته هم‌طول تعریف می‌شود، یعنی تعداد مکانهایی که حاوی حروف نامشابه هستند.

(۲) فاصله لونشتاین^۱ بین دو رشته که الزاماً هم‌طول نیستند، یعنی کمترین تعداد «عملیات ویرایش» مورد نیاز برای تبدیل یک رشته به رشته دیگر، که در اینجا عمل ویرایش عبارت است از حذف، درج، یا تغییر یک حرف در هر توالی. هر توالی داده شده از عملیات ویرایش، یک هم‌ردیف‌سازی یکتا را القا می‌کند. اما عکس این مطلب صحیح نیست.

در زیست‌شناسی مولکولی، می‌دانیم که درج و حذف در توالی‌های زن و پروتئین رخ داده‌اند. بنابراین فاصله همینگ به اندازه کافی عام نیست. به علاوه، شواهدی وجود دارد که نشان می‌دهد موقعیت خود تغییرات محتملت از تغییرات دیگر است. بنابراین حتی فاصله لونشتاین نیز باید تعمیم داده شود تا براساس مدل تکاملی زیربنایی ما، و زنهای متفاوتی برای عملیات مختلف ویرایش در آن منظور گردد. برای مثال، به نظر می‌آید که جهشها محافظه‌کارانه رخ می‌دهند؛ جایگزینی یک آمینواسید در یک پروتئین با آمینواسید دیگری که اندازه یا خواص فیزیکی-شیمیایی یکسان دارد محتملت است تا جایگزینی آن با آمینواسید دیگری که خصوصیات آن بیشتر متفاوت است. برای انکاس این مطلب، به جای شمارش گسسته عملیات ویرایش، به هر تغییر در توالی یک «هزینه» (متعلق به \mathbb{R}) نسبت می‌دهیم. همچنین شواهدی وجود دارد که نشان می‌دهد هزینه یک فاصله خالی، مانند مدل لونشتاین، با طول آن متناسب نیست؛ اگرچه انتخاب مناسب وزنهای فواصل خالی به شکل تابعی از طول دقت زیادی می‌خواهد، در بسیاری از طرحها از یک تابع خطی با یک پارامتر α برای مقدار اولیه دادن به فاصله خالی و یک پارامتر کوچکتر β برای بسط فاصله، به منظور محاسبه هزینه فاصله خالی به شکل

$$(1 - \text{طول فاصله خالی}) \times \alpha + \beta$$

استفاده می‌شود. الگوریتمهای وجود دارند که با کمینه کردن مجموع هزینه‌های عملیات ویرایشی که یک رشته را به رشته دیگر تبدیل می‌کند، بهترین هم‌ردیف‌سازی را تعیین می‌کنند.

مسئله هم‌ردیف‌سازی بهینه توالی را به شکل صوری می‌توان چنین بیان کرد. دو رشته حرفی A و $B = b_1 b_2 \dots b_m$ داده شده‌اند، که هر a_i و b_j عضوی از یک مجموعه الفبای A است. فرض کنید $\{A \cup \emptyset\} = A^+$. یک توالی از عملیات ویرایش، مجموعه‌ای از زوجهای مرتب (x, y) است، که $x, y \in A^+$. عملیات ویرایش متفاوت عبارت‌اند از: جایگزینی b_j با a_i ، که با (a_i, b_j) نشان داده می‌شود، حذف a_i از رشته A ، که با (a_i, ϕ) نشان داده می‌شود، حذف b_j از رشته B ، که با (ϕ, b_j) نشان داده می‌شود. تابع هزینه d بر عملیات ویرایش به شکل زیر تعریف می‌شود:

$$d(a_i, b_j) = \text{هزینه یک جهش}$$

همان نقطه‌ای است که طبیعت موجب می‌شود توالی‌های یک بعدی زنها به جهان سه بعدی‌ای که ما در آن زندگی می‌کنیم جهش کنند.

اهداف زیست‌شناسی مولکولی محاسباتی

اهداف ما از پرداختن به این مبحث چیست؟ نخست آنکه بتوانیم مشابهها و تفاوت‌های موجود بین توالیها و بین ساختارها را به سادگی توصیف و طبقه‌بندی کنیم. تپیلوژی فضای توالی، فضای ساختار، فضای کارکرد پروتئین چیست؟ نگاشتهای بین این فضاهای چگونه‌اند؟ ما می‌خواهیم که با استفاده از تکامل به عنوان اصل سازماندهی، بتوانیم روابط بین توالی، ساختار و کارکرد پروتئین را توصیف و پیشگویی کنیم.

از کجا باید کار را شروع کرد؟ سیدنی برز یک بار با افسوس گفت: «مسئله در زیست‌شناسی این است که هیچ نوسانگر همسازی وجود ندارد.» منظور او از این گفته این بود که در زیست‌شناسی، برخلاف فیزیک، راه گریزی از پیچیدگی نیست، حتی از طریق ایده آل سازی. در فیزیک، نوسانگر همساز یک مسئله ساده است که می‌توان به روشهای متعددی آن را به دقت حل کرد؛ این مسئله در مورد برخی از پدیده‌ها دقیقاً قابل اعمال است، و برای سایر پدیده‌ها تخمین مؤثری است. نوسانگر همساز در فیزیک یک بستر آزمون سنتی برای روشهای جدید است. در حقیقت، زیست‌شناسی مولکولی محاسباتی دو «نوسانگر همساز» به تعبیر برز دارد: هم‌ردیف‌سازی توالیها، و برهمنهی ساختارها. این اعمال، که می‌توان آنها را با دقت و کارایی انجام داد، پایه بسیاری از تحلیلهای روابط توالی-ساختار در زیست‌شناسی مولکولی هستند. اکنون دانستن این نکته که جهان واقعی اغلب ناهم‌ساز است شگفت‌انگیز نخواهد بود. با این‌همه، بسیاری از این ابزارهای ارزشمند بدکمک این موارد ساده ساخته شده‌اند.

اما، ابزارها پاسخ می‌سازند نه سؤال. پژوهش در این حوزه همچنان بر تعامل بین داشمندان و داده‌ها، با کمک روشهای ریاضیاتی و محاسباتی، متکی است.

توالیها و هم‌ردیف‌سازی آنها

توالی‌های زن و پروتئین شکل رشته‌های حرفی به‌خود می‌گیرند. برای توالی‌های زن، حروف از مجموعه چهار عضوی $\{A, T, G, C\}$ انتخاب می‌شوند که نشان‌دهنده نوکلئوتیدهای آدنین، تایین، گوانین و کایتوزین هستند. برای توالی‌های پروتئین، حروف از یک مجموعه بیست عضوی انتخاب می‌شوند که نشان‌دهنده بیست آمینواسید معمولی هستند.

هم‌ردیف‌سازی

هم‌ردیف‌سازی دو رشته حرفی به معنی تعیین یک تاظر معنی دار بین مؤلفه‌های آنهاست. برای دو رشته حرفی زیر

c t a t a a t c g c t g a a c

دو هم‌ردیف‌سازی ممکن عبارت‌اند از

g c t g - a a - c g c t g a - a - - c

و

- c t a t a a t c - c t - a t a a t c

TYLWEFLLKLLQDR. EYCPRFIKWTNREKGVFKLV.. DSKAVSRLWGMHKN. KPD VQLWQFLLILEILTD.. CEHTDVIEWVG. TEGEFKLT.. DPDRVARLGEEKNN. KPA IQLWQFLLLELLTD.. KDARDCISWVG. DEGEFKLN.. QPELVAQKWCQRKN. KPT IQLWQFLLLELLSD.. SSNSCITWEG. TNGEFKMT.. DPDEVARWRGERKS. KPN IQLWQFLLLELLTD.. KSCQSFISWTG. DGWEFKLS.. DPDEVARRWGKRKN. KPK IQLWQFLLLELLQD.. GARSSCIRWTG. NSREFQLC.. DPKEVARLGCRK. KPG IQLWHFILELLQK.. EEFRHVIAWQQGEYGEFIK. DPDEVARLGRRK. KPQ VTLWQFLLQLLRE.. QGNHIIISWTSRDOGEGFKLV.. DAEELVARLGRLKN. KTN ITLWQFLLHLLD.. QKHEHHLICWTS. NDGEFKLL.. KAEELVALWGLRKN. KTN LQLWQFLLVALLD.. PTNAHFIATWG. RGMEFKLI.. EPEEVARLGQIJKN. RPA IHLWQFLLKELLASP. QVNCTAIRWIDRSKGIFKIE.. DSVRAVAKLWGRKN. RPA RLLWDFLQQLNDRNQKYSQSLIAWKCRDTGVFKIV.. DPAGLAKLGQIJKN. HLS RLLWDYVYQLLSD.. SRYENFIRWEDKESKIFRIV.. DPNGLARLWGNHKN. RTN IRLYQFLLDLLRS.. GDMKDSIWWVVDKDKGTQFQSSKKHEALAHRGQIJKGNRK LRLYQFLLGLTR.. GDMRECWWVVEPGAGVQFQSSKKHELLARRWGQQKGQRK

L f1 1L i W F a WG K

شکل ۲ هم ردیف‌سازی چندگانه توالیهای ناتمام یک خانواده از پروتئینها که دامنه‌های ETS نامیده می‌شوند. هر خط با توالی آمینواسیدهای یک پروتئین متناظر است، که با توالی ای از حروفی که هر کدام نماینده یک آمینواسید است مشخص می‌شود. با نگاه به هر ستون می‌توان آمینواسیدهایی را که در آن مکان در هر یک از پروتئینها خانواده پدیدار می‌شوند مشاهده کرد. به این طریق الگوهای اولویت‌دار قابل رویت می‌شوند. مثلاً مکان سوم حاوی یک لوسین یا I در هر توالی است – این امر دلالت برین دارد که برخی محدودیتهای ساختاری یا کارکردی، تکامل را از تغییر وضعیت این مکان بازداشت‌اند. حروف زیر جدول مکان سازه‌های ناتماین (حروف بزرگ) یا ناتماین با یک استنتا (حروف کوچک) را مختص می‌کنند. به توزیع نابرابر تغییر در سنتهای مختلف توجه کنید. تابو سازه‌های ماندگار ۴، ۳، ۲، ۱ و ۰ از وجود مارپیچها در پروتئین خبر می‌دهد، که خبر صحیح است. الگوهای دیگر عیقطر پنهان شده‌اند، و تعین آنها نیاز به تحلیل محاسباتی دارد. چنین الگوهایی ممکن است حاوی همیستگیهایی میان توزیع آمینواسیدها در مکانهای مختلف باشند. مثلاً در ستون چهارم از سمت چپ آمینواسید تیروزین یا Y، تنها در دو توالی آخر پدیدار می‌شود، سایر توالیها دارای تریوتوفان یا W هستند. یک همبستگی تقریبی در الگوی تغییر بین این ستون و سنتهای چهارم و پنجم از راست وجود دارد. بسیاری عقیده دارند (با حداقل امیدوارند) که همیستگیهای الگوهای تغییر در مکانهای مختلف چنین جدولی از توالیها، سرنجهایی در باره محلهایی که در فضای سه‌بعدی بر هم اثر می‌گذارند به ما بددهد. متأسفانه قوانن بسیار ضعیف است.

خواص فیزیکی-شیمیایی یکسان می‌شود – نتایج ساختاری تغییرات توالی را تعديل می‌کند.

حتی اگر مشابهت در سطح توالی قابل تشخیص باشد، برای پروتئینها که ارتباط دوری با هم دارند، به کمک مقایسه ساختارها که آخرین راه چاره است، می‌توان فهمید که هم ردیف‌سازی جفت جفت بهینه توالیها اغلب منجر به جواب غلط می‌شود.

اما اگر تعداد زیادی توالی مرتبط در دسترس باشد، هم ردیف‌سازی چندگانه توالیها نتایج با ارزشتر و دقیقتری نسبت به هم ردیف‌سازی جفت جفت توالیها به دست می‌دهد. چرا هم ردیف‌سازیهای چندگانه اطلاعات توالی را گسترش می‌بخشند؟ از اینجا الگوهای ماندگاری ظاهر می‌شوند. گستره و طبیعت تغییر در هر یک از مکانها راهنمای مهمی برای تعیین نقش ساختاری یا کارکرد ناحیه‌های مختلف توالی است (شکل ۲). برای مثال، سازه^۱هایی که در یک خانواده کامل از پروتئینها ماندگار [=ب[Tغییر] بوده‌اند معمولاً در کارکرد دخیل‌اند، یا دستکم اغلب نقشی اساسی در ساختار ایفا می‌کنند. بر عکس، مناطقی که عملیات درج و حذف در آنها زیاد صورت می‌گیرد معمولاً با

$d(\emptyset, b_j) = \text{هزینه یک حذف یا درج}$

و کمترین فاصله وزن دار بین رشته‌های A و B عبارت است از

$$D(A, B) = \min_{A \rightarrow B} \sum d(x, y)$$

که در آن $x, y \in A^+$ و مقدار مینیمم بدارای همه توالیهای عملیات ویرایش که A را به B تبدیل می‌کنند محاسبه می‌شود. اگر (y) یک متریک بر A^+ باشد، $D(A, B)$ متریکی بر رشته‌های حروف از A^+ خواهد بود. (در این بیان از مسأله فرض می‌شود که هزینه فاصله‌های خالی مستقل از طول آنهاست؛ طرحهای واقع‌بینانه‌تر که به فاصله‌های خالی وزن می‌دهند، تعییمی از این طرح هستند).

مسأله یافتن $D(A, B)$ است و یک یا چند هم ردیف‌سازی که با آن متناظرند. الگوریتمی که این مسأله را در زمان $O(mn)$ حل می‌کند زمان درازی است که شناخته شده است، و در سیاری از مسائل از قبیل ویرایش متن، تشخیص گفتار، و تحلیل آواز یوندگان به کار گرفته شده است [۲]. این الگوریتم توسط مقاله تأثیرگذار نیدلمن و وونش [۳]، به زیست‌شناسی وارد شد. چند خصوصیت این الگوریتم قابل توجه‌اند.

- این الگوریتم یک بهینه مطلق به دست می‌دهد: توجه داشته باشید که این یکی از دو «نوسانگر همساز» زیست‌شناسی مولکولی محاسباتی است. ما رویی در اختیار داریم که مطمئنیم که در مینیمهای موضعی به دام نخواهد افتاد.
- این خبر خوب بود. خبر بد اینکه تعییر نتایج چندان سراسرت است. اگرچه توالی ای از عملیات ویرایش که از یک هم ردیف‌سازی بهینه ناشی شده است ممکن است با یک مسیر تکاملی واقعی متناظر باشد، اما اثبات اینکه چنین است ممکن نیست. هر چه فاصله ویرایش بیشتر باشد، تعداد مسیرهای تکاملی معقول بیشتر است. نه تنها هم ردیف‌سازی‌های بهینه ممکن است یکتا نباشند، بلکه ممکن است هم ردیف‌سازیهای زیربهینه بسیاری وجود داشته باشند که ارزش آنها کاملاً به مقدار بهینه نزدیک باشد. مثلاً فیج و اسیمیت ژنهای جوجه را برای هموگلوبینهای α و β آزمایش کردند [۴]. آنها ۱۷ هم ردیف‌سازی بهینه یافتند، که یکی از آنها با هم ردیف‌سازی مبتنی بر ساختارهای هموگلوبین شناخته شده مطابقت داشت، و بیش از هزار هم ردیف‌سازی یافت شد که با مقدار بهینه کمتر از ۵٪ اختلاف داشت.

مشکلات ناشی از هم ردیف‌سازی جفت جفت توالیها
مشاهده شده است که با تکامل پروتئینها، توالیهای آمینواسیدها بسیار سریعتر از ساختار واگرا می‌شوند. در بسیاری از موارد می‌توان یک رابطه تکاملی بین دو ساختار پروتئین یافت، حتی اگر هیچ شباهت قابل درکی بین توالیهای زنها یا توالیهای آمینواسیدها وجود نداشته باشد. آنچه اتفاق می‌افتد این است: زنها فضای توالیهای DNA را می‌کاوند، اما انتخاب طبیعی به مثابة ترمیزی در برابر تغییر ساختار عمل می‌کند تا کارکرد را حفظ کند. افزونگی^۱ در کد زنیک – این واقعیت که تعداد زیادی از بازهای سه‌گانه یک آمینواسید یکسان را کدنویسی می‌کنند، و بسیاری از تغییرات تک بازی منجر به حصول آمینواسیدهایی با

1. redundancy

همه ژنها را در اختیار ما می‌گذارد. تنها برای اقیت کوچکی از این ژنها ساختار سه بعدی پروتئینهای متناظر را در اختیار داریم.

تحلیل ساختار پروتئینها

اولین مشکل تحلیل ساختار مولکولهایی به پیچیدگی پروتئینها، مشکل بازنمایی است. تکنیکهای گرافیک کامپیوتری برای رسم بازنمایهای ساده شده‌ای از پروتئینها ارائه شده‌اند. شکل ۳ روش می‌کند که برای یک مولکول کوچک پروتئین، تعبیر یک بازنمایی دقیق با جزئیات کامل چقدر دشوار است، و نوع تصاویر ساده شده‌ای را که برنامه‌ها تولید می‌کنند تا امکان دستیابی بصری به موضوع را به ما بدهند نشان می‌دهد. آزمایشگاههای کوچک متعدد، تعداد زیادی بازنمایی متفاوت تولید کرده‌اند؛ یعنی افراد بسیاری بازنمایهای ساده شده مختلفی را پیشنهاد کرده‌اند، و این پیشنهادات در بسته‌های گرافیکی عالی جمع‌آوری شده است. یک ترسیمگر مولکول ماهر آنها را با هم ترکیب خواهد کرد تا جنبه‌های مختلف یک ساختار را با درجات دقت قابل تنظیم نمایش دهد. چنین تصاویری، که به صورت تمام رنگی درآمده‌اند و در آنها از جلوه‌های سایه‌زنی فانتزی (اما غیرواقعی)، با توجه به اندازه مولکولها نسبت به طول موج نور مرئی استفاده شده است، مجلات، بوسترها، حتی لباسها و فنجانها را می‌آرایند. ما اکنون ساختار ۱۰۰۰۰ پروتئین را می‌شناسیم، و طیف وسیعی از الگوهای فضایی را در آنها مشاهده می‌کنیم. در پاسخ به حرف راترفورد که معتقد است «همه علم یا فیزیک است یا جمع کردن تمیز»، من می‌گویم که مطالعه ساختار پروتئینها بهترین جنبه‌های این دو رشته را در هم ادغام کرده است! ما با تنویر چشمگیری مواجه‌هیم، اما در عین حال به وجود اصول جامع در پس زمینه ایمان داریم.

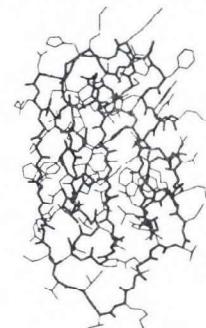
هر پروتئین از یک زنجیره اصلی پلیمر تکرارشونده خطی (یعنی بدون شاخه) تشکیل شده است که شاخه‌های جنبی آمینواسید در فواصل منظم به آن آویخته‌اند. پروتئین با رشته‌ای از چراغهای درخت کریسمس مقایسه است، که سیم آن با زنجیره اصلی تکرارشونده و رشته‌های رنگی چراغها با شاخه‌های جنبی متفوّغ متنوّع متناظرند.

زنجیره اصلی یک خم فضایی را توصیف می‌کند که به وسیله تعاملهای مطلوب میان زنجیره‌های جنبی که به هم متصل شده‌اند پایدار شده است. چنین خم فضایی را در بخش میانی شکل ۳ برآحتی می‌توان دید. دو منطقه در جلو تصویر فرم ماریچی دارند که محورشان تقریباً عمودی است. این یکی از دو ساختار استاندارد دیگر، نوار تقریباً بازشده است: پروتئین شکل ۳ چهار رشته نوار در خود دارد که در راستای تقریباً عمودی قرار دارند. این نوارها از پهلو بر هم اثر می‌گذارند تا اجتماع خود را پایدار کنند. در چارچوب پایینی شکل ۳، ماریچها و نوارها، با «نماد»‌هایی نشان داده شده‌اند، ماریچها با استوانه و نوارها با فالشهای بزرگ. چارچوب بالایی شکل ۳ بازنمایی ساختار را با بیشترین جزئیات، شامل زنجیره اصلی و شاخه‌های جنبی نشان می‌دهد؛ تضاد رنگها نشان‌دهنده اهمیت ساده‌سازی در تولید تصویر حتی از یک پروتئین کوچک است، بهنحوی که تصویر به لحاظ بصری قابل فهم باشد.

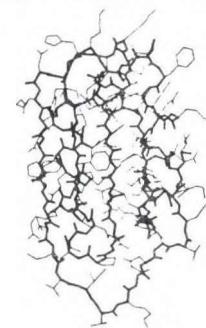
گام اولیه تجزیه یک ساختار جدید عبارت است از تعیین مناطق ماریچ و نوار. این اطلاعاتی است که برای تبدیل بازنمایی چارچوب مرکزی شکل ۳

مناطق حاشیه‌ای متناظرند. (یک توالي دیگر ساختار خود متفعل است، دل جفت توالي هم دفعه ساختار خود را نجوا می‌کنند، سه با جند توالي ساختار خود را به صدای بلند فریاد می‌کنند).

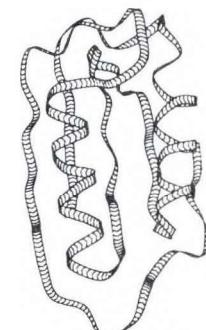
اگر تواليها تنها به صورت غیرمستقیم به ساختار اشاره می‌کنند، چرا مستقیماً به سراغ ساختار نرویم؟ علت این است که مقدار داده‌ها در مورد تواليها شناخته شده بسیار بیشتر از مقدار داده‌های ساختاری است. برای حدود ۲۰ موجود زنده، توالي همه ژن‌های تعیین شده است، که تواليها کامل



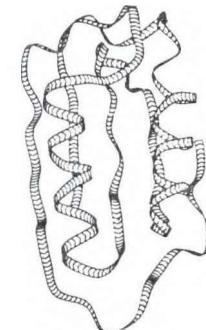
آسیل فسفات



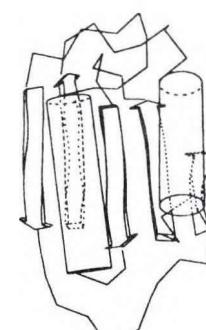
آسیل فسفات



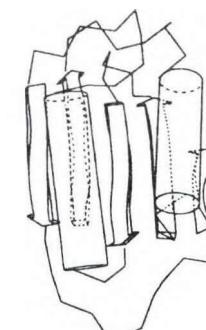
آسیل فسفات



آسیل فسفات



آسیل فسفات



آسیل فسفات

شکل ۳ پروتئینها چنان ساختارهای پیچیده‌ای هستند که لازم آمده است تا ابراهای ویژه‌ای برای نمایش آنها ساخته شود. این شکل یک پروتئین نسبتاً کوچک به نام آسیل فسفات را با سه درجه مختلف ساده‌سازی نشان می‌دهد. بالا: مدل کلی کامل؛ زنجیره اصلی پرنگک از شاخه‌های جنبی است. وسط: مسیر زنجیر، توسط یک خم درون‌بایی شده هموار بازنمایی شده است، فلشها هشت زنجیره را تعیین می‌کنند. پایین: نمودار اجمالی، که در آن استوانه‌ها نمایانگر ماریچها و فلشها نمایانگر نوارهاستند. با تغییر خطوطی که از پشت اشیاء صلب می‌گذرند به خطوط شکسته، این اشیاء به صورت «نیمه‌شفاف» بازنمایی شده‌اند. برای برهنه‌ی بازنمایهای مجاور، صفحه را ۹۰° بچرخانید و به صورت سه بعدی نگاه کنید (ولی نه خیلی طولانی!).

مسائل (۲) و (۳) نیازمند تعیین نحوه هم‌ردیف‌سازی نقاط هستند. روش‌های هم‌ردیف‌سازی که منحصراً بر مختصات (نه بر توالی‌های آمینو اسید) مبتنی هستند، هم‌ردیف‌سازی ساختاری نامیده می‌شوند. در هم‌ردیف‌سازی‌های ساختاری، سازه‌های متناظر یکی گرفته می‌شوند زیرا آنها نسبت به کل ساختار مکان مشابهی اختیار می‌کنند. باید به استخراج زیرساختار مشترک ماکسیمال و بنا نهادن هم‌ردیف‌سازی بر آن اندیشید. (برای مثال، زیرساختار مشترک ماکسیمال حروف B و R، حرف P است). سازه‌های خارج از زیرساختار مشترک ماکسیمال قابل هم‌ردیف شدن نیستند، حقیقتی که با هم‌ردیف‌سازی

جفت‌جفت توالیها مشخص نمی‌شود؛ این یکی از ضعفهای آن است. کلیترین رویکرد به این سه مسئله بر حل مسئله (۱)، یعنی حالت تناظر معلوم $p_i \leftrightarrow q_i$ مبتنی است. دو شیء یکسان را می‌توان با انتقال و دوران صلب یکی از آنها، برهم‌نهاد. دو شیء که هشایه هستند از طریق دوران و انتقال به برهم‌نهی تقدیب می‌رسند. اگر اشیاء مجموعه‌های مرتب نقاط باشند، شاخص مشابهت آنها برابر با جذر میانگین مربع انحرافها، Δ ، پس از برهم‌نهی بهینه خواهد بود:

$$\Delta^* = \min_{R, t} \left\{ \sum_{i=1}^N \| Rp_i + t - q_i \|^2 \right\}$$

که در آن R ماتریس دوران مناسب ($\det R = 1$) و t بردار انتقال است. در برهم‌نهی بهینه، مکانهای میانی (به زبان محاوره، «مراکز ثقل») دو مجموعه برهم منطبق می‌شوند. مسئله تعیین جهت نسبی صحیح به عنوان «مسئله پروکروتس معتمد» شناخته شده است و راه حل‌هایی مبتنی بر روش‌های استاندارد جبر خطی برای آن وجود دارد [۵].

حل مسئله زیرساختار مشترک ماکسیمال مبنای تعریف یک متريک را برای ساختارها فراهم می‌کند. بر اين مبنای توان مشابههای مقطعی و جزئی را پیدا کرد، و يك درخت طبقه‌بندی برای كل مجموعه ساختارهای پروتئین ارائه داد. رویکردهای موجود به محاسبه زیرساختار مشترک ماکسیمال بر دو بازنمایی از ساختارها متکی بوده‌اند: (۱) به صورت فهرستهایی از مختصات $p_i = (x_i, y_i, z_i)$ ، $i = 1, \dots, n$ ، (۲) به صورت ماتریسهای فاصله $|p_i - p_j| = D_p(i, j)$. مزیت عمده ماتریسهای فاصله این است که يك بازنمایی مستقل از مبدأ و راستا برای ساختار ارائه می‌کنند. بر حسب ماتریسهای فاصله، مؤلفه ماکسیمال تقاضل ماتریسهای فاصله $D_q(i, j) = \max_{1 \leq i < j \leq n} |D_p(i, j)|$ ، شاخصی برای اختلاف ساختاری بین دو مجموعه نقطه هم‌ردیف شده به دست می‌دهد.

مختصات و ماتریسهای فاصله، بازنمایهای تقریباً معادل از یک مجموعه نقاط هستند. محاسبه ماتریس فاصله از روی مختصات بدیهی است. اینکه آیا می‌توان مختصات را دقیقاً و مستقیماً از روی ماتریس فاصله بازسازی کرد کمتر واضح است، اما این کار را می‌توان به کمک قطعی سازی ماتریس انجام داد [۶]. البته، ماتریس فاصله هم ساختار اولیه و هم تصویر آینه‌ای آن را به طور معادل مشخص می‌کند (پس دسته‌های چپ و راست متضاد، دو تصویر آینه‌ای هستند)، اما این ابهام مشکلی جدی برای کاربردهای زیست‌شناسی مولکولی نیست. اطلاعات مربوط به موقعیت و جهت هم مسلماً مفقود می‌شوند.

به چارچوب پایینی مورد نیاز است. متدالترین نوع مارپیچ در ساختارهای پروتئینی، در هر پیچ ۳۶ سازه را دربرمی‌گیرد. خصوصیاتی از توالی که این تناوب را نشان می‌دهند، مناطق مارپیچی را تداعی می‌کنند.

برهم‌نهی ساختارها

مانند توالیها، مسئله اساسی در تحلیل ساختارها هم طرح و محاسبه یک شاخص مشابه است. فرض کنید که مجموعه‌هایی مختصاتی در اختیار داریم که دو ساختار را بازنمایی می‌کنند:

$$p_i = (x_i, y_i, z_i), \quad i = 1, \dots, N$$

$$q_j = (x'_j, y'_j, z'_j), \quad j = 1, \dots, M$$

درست مانند توالیها، در اینجا هم مسئله هم‌ردیف‌سازی مطرح می‌شود. تقابل بین سه مسئله مرتبط را که در شیمی محاسباتی مطرح آند در نظر بگیرید.

(۱) شاخص مشابهت دو مجموعه از اتمها با تناظرهای داده شده

$$p_i \leftrightarrow q_i \quad i = 1, \dots, N$$

(۱) تعیین کنید (همتای این شاخص برای توالیها، فاصله همینگ است). این مسئله را می‌توان دقیقاً و به نحو کارآیی حل کرد — این دومین «نوسانگر همساز» زیست‌شناسی مولکولی محاسباتی است.

(۲) شاخص مشابهت دو مجموعه از اتمها با تناظر نااعلوم (۱) تعیین کنید، اما با این فرض که ساختار مولکولی آنها — به طور خاص، ترتیب خطی سازه‌ها — تناظر را محدود می‌کند. در مورد پروتئینها، هم‌ردیف‌سازی باید ترتیب را در طول زنجیره حفظ کند:

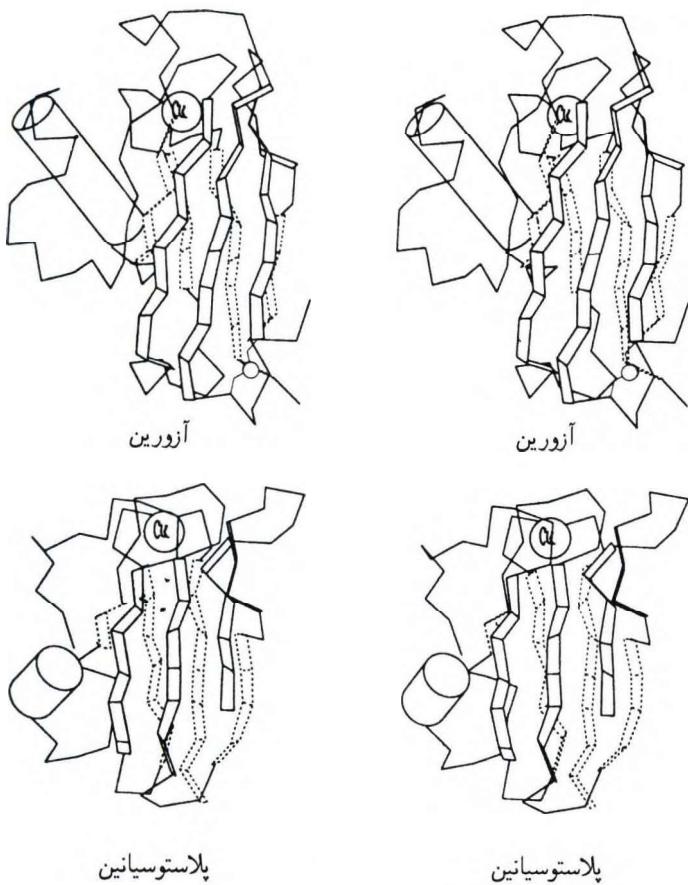
$$p_{i(k)} \leftrightarrow q_{j(k)}, \quad k = 1, \dots, K \leq N, M$$

با این قید که $i(k_1) > i(k_2) \Rightarrow k_1 > k_2$ و $j(k_1) > j(k_2)$ ، این شاخص را می‌توان متضاد با فاصله لونشتاین بین رشته‌های حرفی، یا هم‌ردیف‌سازی توالیها با فواصل خالی در نظر گرفت.

(۳) شاخص مشابهت بین دو مجموعه از اتمها با تناظرهای نااعلوم، بدون هیچ محدودیتی، بر تناظر، (۱) تعیین کنید:

$$p_{i(k)} \leftrightarrow q_{j(k)}, \quad k = 1, \dots, K \leq N, M$$

این مسئله در حالت مهم زیر پیش می‌آید: فرض کنید دو (یا چند) مولکول تأثیرات زیست‌شناسختی مشابه دارند، مثلاً دارای خواص دارو-ساختاری یکسان هستند؛ این حالت معمولاً وقتی پیش می‌آید که ساختارها در زیرمجموعه‌ای نسبتاً کوچک از اتمهایشان که مسؤول فعالیت زیستی معینی است، و یک فازماکوفود نامیده می‌شود، مشترک باشند. برای تشخیص فارماکوفور، خوب است بتوانیم زیرمجموعه‌های ماکسیمال اتمهای دو یا چند مولکول را که ساختار مشابه دارند بیابیم.



شکل ۴ در حین تکامل، جهش‌های توالیهای Ζن بر روی هم انباشته می‌شوند، و در نتیجه توالیها و ساختارهای پروتئین از هم فاصله می‌گیرند. این شکل دو پروتئین مرتبط انتقال هندۀ الکترون، یعنی پلاستوسیانین برگ سپیدار و آزورین باکتریایی را نشان می‌دهد. بخشی از ساختارها که در دینامیک راست تصویر واقع است و شامل آرایه‌های پرو-ونتھچین مناطق «روبانی شکل» — که نوار نامیده می‌شوند — و محل پیوند با من است، به خوبی در حین تکامل ماندگار است، در حالی که بخشی از ساختار که در دینامیک چب تصویر واقع است به شدت تغییر کرده است.

آنکه توالیهای آمینواسید انسولین خوک و انسان یکسان نیستند. اعتماد به چنین مشابههایی روشی برای پیش‌بینی ساختار پروتئینها از روی ساختارهای بسیار نزدیک شناخته شده به دست می‌دهد، که به عنوان «مدلسازی مانستگی»^۱ شناخته شده است. اما، همچنان که تکامل ادامه می‌یابد، توالیها و ساختارها نهایتاً به شکل بنیادیتری واگرای می‌شوند. شکل ۴ دو پروتئین پلاستوسیانین و آزورین را نشان می‌دهد که با هم ارتباط کمی دارند، در منطقه راست شکل، دو نوار وجود دارد که روبرو بسته شده‌اند، و «هسته» ماندگار ساختار را تشکیل می‌دهند، در حالی که منطقه حلقی بلند سمت چپ از شکلی کاملاً متفاوت در دو ساختار برخوردار است.

پیش‌بینی ساختار پروتئین

طیعت الگوریتمی در اختیار دارد که به کمک آن ساختار سه‌بعدی پروتئین را فقط از روی توالی آمینواسید آن تعیین می‌کند. ما باید قادر به کشف آن

دشواری اصلی در محاسبات زیرساختار مشترک ماکسیمال از نوع (۲) و (۳)، پیچیدگی ترکیباتی ناشی از در نظر گرفتن همه هم‌ردیف‌سازیهای ممکن است. معلوم شده است که الگوریتمهای مبتنی بر ماتریسهای فاصله در مواجهه با این مسئله بسیار کارآمدتر از الگوریتمهای مبتنی بر مجموعه‌های مختصات عمل می‌کنند. بازنمایهای ماتریسی مرتبط که به جای مختصات اتمی، بر مؤلفه‌های ساختاری نظری مارپیچها و نوارها مبتنی هستند، بازنمایی فشرده‌ای از الگوهای تاخورده‌گی پروتئینها به دست می‌دهند. استخراج زیرماتریسهای مشترک ماکسیمال، بزرگترین زیرساختارهای دارای الگوی تاخورده‌گی مشترک را آشکار می‌کند. به علاوه چنین بازنمایهایی، امکان شمارش همه الگوهای تاخورده‌گی پروتئین را در اختیار ما می‌گذارد. به صورت تجربی برآورد شده است که همه پروتئینهای طبیعی کمتر از حدود ۱۰۰۰ الگوی تاخورده‌گی دارند. شمارش کامل به ما امکان می‌دهد که انتخابهای طبیعت را آزمایش کنیم، و سعی کنیم تصادف تاریخی و ضرورت معماری را از هم تمیز دهیم.

تکامل پروتئین

مطالعه تکامل پروتئین یعنی بررسی اینکه چگونه توالیهای آمینواسید و ساختارهای پروتئین متناظر در گونه‌های وابسته با هم اختلاف پیدا می‌کنند. این تحقیقی از نوع اطلاع‌دهنده است، که به ما در فهم روابط توالی-ساختار کمک می‌کند. زیرا اگرچه ما می‌دانیم که یک توالی آمینواسید متفاوت همه اطلاعات لازم برای مشخص کردن ساختار پروتئین را در خود دارد، هنوز نمی‌فهمیم که چگونه باید از توالی به ساختار رسید. این را «صورت انتگرالی» مسئله تاخورده‌گی پروتئین در نظر بگیرید، که مسئله حل نشده‌ای است. در بررسی تکامل پروتئین، متناهده می‌کنیم که چگونه تغییرات توالیها در تغییرات ساختار منعکس می‌شوند؛ این را آسانتر می‌توان فهمید. این مسئله را «صورت دیفرانسیلی» مسئله تاخورده‌گی پروتئین در نظر بگیرید.

موضوع:	تاخورده‌گی پروتئین	تکامل پروتئین
مساهده:	توالی ← ساختار	تغییر در توالی ← تغییر در ساختار
حالت مسئله:	«صورت انتگرالی»	«صورت دیفرانسیلی»
وضعیت مسئله:	حل نشده	حل نشده اما جاید ساده‌تر باشد

استدلال ساده‌ای بر این دلالت می‌کند که ساختار باید تابعی تقریباً «پیوسته» از توالی باشد، حداقل برای توالیها و ساختارهایی که به صورت طبیعی شکل گرفته‌اند. فرض کنید پروتئینی وجود می‌داشت که در آن هر جهشی (هر تغییری در توالی آمینواسید) یک ساختار ناپایدار تولید می‌کرد. در این صورت طبیعت هرگز نمی‌توانسته است با فرآیندهای تکاملی به چنین ساختاری برسد، زیرا هیچ صورت اولیه پایداری نمی‌توانسته برای آن وجود داشته باشد. پس چنین نتیجه می‌شود که ساختارهای طبیعی باید مستحکم باشند. اغلب تغییرات کوچک در توالی باید تغییری در ساختار ندهند. (این شرط برای ساختارهای پروتئینی که به روش‌های مصنوعی ساخته می‌شوند برقرار نیست.) در واقع، پروتئینهای طبیعی با توالیهای بسیار شبیه به هم، ساختارهای بسیار شبیه به هم دارند. پیش از آنکه انسولین مصنوعی انسانی در دسترس قرار گیرد، انسولین خوک درمان بالینی مؤثری برای انسانهای مبتلا به مرض قند بود، با

درواقع شما چگونه می‌توانید مردم را مقاعده کنید که روش موقتی برای پیش‌بینی ساختار پروتئین در اختیار دارید؟ دو نوع از ادعاهای از اساس غیرقابل آزمون هستند. یکی آنکه شما می‌توانید ساختار پروتئینی را که ساختار آن از قبل دانسته شده است پیش‌بینی کنید. دیگر آنکه شما ساختار پروتئینی را پیش‌بینی کرده‌اید که ساختار تجربی آن ناشناخته است و به احتمال زیاد تا مدت طولانی ناشناخته خواهد ماند. باید در حوزه بینایی میان شناخته‌شده‌ها و آنچه تا زمان طولانی غیرقابل شناسایی است کار کرد، و پیش‌بینی ساختار را با تعیین ساختارهای در دست اقدام هماهنگ نمود.

برای نظم بخشنیدن به این فعالیت، برای پاداش دادن به کسانی که موجب پیشرفت‌های اصیل شده‌اند، و برای تکذیب ادعاهای آنانی که مصراً ادعا می‌کرند «مسئله پیش‌بینی ساختار پروتئین را حل کرده‌اند»، جان مالت^۱ ایده آزمونهای کور سازمان یافته را مطرح کرد. ایده این است که دانشمندان در حین فرایند کشف ساختارها، توالیهای آمینواسید را به همه اعلام کنند، اما قول بدھند که ساختار را تا زمان مورد توافقی مخفی نگه دارند. همه کسانی که باور دارند که روشی برای پیش‌بینی ساختار پروتئین در اختیار دارند می‌توانند پیش‌بینیهای خود را تا قبل از تاریخ اعلام ساختار ارائه دهند. پس از اعلام، می‌توان پیش‌بینیها را با تجربه مقایسه کرد — که این باعث خرسنیدی عده اندکی و تأسف گروه کثیری خواهد شد. این ایده تحت برنامه CASP — برآورد انتقادی پیش‌بینی ساختار^۲ — براساس یک دوره دو ساله ارائه شده است.

روشهای پیش‌بینی به دو طبقه کلی تقسیم می‌شوند، استقرایی و استنتاجی. در روشهای استقرایی مستقیماً از بانکهای داده توالیها و ساختارها استفاده می‌شود. روشهای استنتاجی، رویکردهای واقعاً ابتدا به ساکن هستند — حالت جزیره‌بی‌آب و علف — که هدف آنها پیش‌بینی ساختار پروتئین براساس اصول عام فیزیک، شیمی و زیست‌شناسی، بدون ارجاع صریح به توالیها و ساختارهای شناخته‌شده است. البته پیشرفت روشهای ابتدا به ساکن بستگی به چیزهایی دارد که از بررسی توالیها و ساختارها آموخته‌ایم. تایز این روشهای در این است که درک حاصل از این بررسی به شکل اصول عامی خلاصه شده است که می‌توان آنها را بدون مراجعه به اطلاعات خاص در پایگاه‌های داده، بدکار برد.

روشهای پیش‌بینی ابتدا به ساکن را می‌توان به دو رویکرد تقسیم کرد، که من آنها را «طبیعی» و «зорی» می‌نامم. رویکرد طبیعی به دنبال درک فرایند تاخورده‌گی طبیعی می‌گردد و سپس سعی می‌کند آن را شبیه‌سازی کند. رویکرد «зорی» هر فرایندی را که بتواند زنجیره را به شکل [=ترکیب، ساختمان] مناسب برساند مجاز می‌داند، حتی اگر این فرایند در راستای مسیری انجام شود که طبیعی نباشد یا حتی از نظر فیزیکی غیرقابل تحقق باشد.

شواهدی هست که انتخاب طبیعی نه تنها وضعیت طبیعی نهایی پروتئینها، بلکه مسیر تاخورده‌گی آنها را نیز شکل داده است. زیرا نه تنها پروتئینها باید به گونه‌ای تکامل یافته باشند که شکل فعل پایداری پیدا کرده باشند، بلکه باید در زمان معقولی، از یک وضعیت تاخورده اولیه شامل مخلوطی از شکلهای تصادفی به چنین شکلی دست یافته باشند. یک محاسبه ساده براساس

باشیم. سپس باید بتوانیم ساختار پروتئینها را که در توالیهای ژن انسان و سایر ژنومها به صورت ذاتی قرار دارد پیش‌بینی کنیم، و آنها را در مسائل عملی نظیر طراحی دارو به کار گیریم. معلوم شده است که مسئله پیش‌بینی ساختار پروتئین مسأله بسیار دشواری است. رویکردهای بسیاری اختیار شده‌اند، و ادعاهای بسیاری مطرح شده‌اند. اما در حال حاضر هیچ روش محاسباتی که بتواند به صورت سازگار حتی یک پیش‌گویی صحیح کیفی از ساختار پروتئین براساس توالی آمینواسید ارائه کند وجود ندارد، مگر آنکه یک پروتئین بسیار مشابه موجود باشد.

فرض کنید که توالی آمینواسیدهای یک پروتئین جدید به شما داده شده بود، و از شما خواسته بودند که ساختار آن را پیش‌بینی کنید. شما باید سعی می‌کردید چه چیزی را پیش‌بینی کنید؛ کاملترین اطلاعاتی که یک پیش‌بینی ممکن است به دست دهد مجموعه کاملی از مختصات سه‌بعدی مدل یک پروتئین نهایی است — یعنی یک پیش‌بینی سه‌بعدی. یک هدف کمتر بلندپروازانه می‌تواند پیش‌بینی این باشد که مناطق مارپیچ و نوار در کجا توالی پدیدار می‌شوند — یعنی یک پیش‌بینی یک بعدی. جایی بین این دو، پیش‌بینی‌هایی قرار دارند که از پیش‌بینیهای ساختار ثانویه یک بعدی فراتر می‌روند، اما تنها برخی اطلاعات کیفی در باره آرایش فضایی کلی الگوی تاخورده‌گی ارائه می‌کنند — باید اینها را پیش‌بینیهای دو بعدی بخوانیم.

از چه نوع اطلاعاتی می‌توان در پیش‌بینی ساختار پروتئین استفاده کرد؟ هدف نهایی، رویکرد ابتدا به ساکن^۳ «محض» است — تنها از توالی پروتئین مقصد استفاده کنید و نه هیچ چیز دیگر. باید به خاطر داشت که این همان کاری است که طبیعت می‌کند — پروتئینها وقتی می‌خواهند تا شوند در پایگاه‌های داده و ب جستجو نمی‌کنند. اما ما می‌توانیم جستجو کنیم، و موقفيت‌هایی هم در استفاده از اطلاعات بانکهای داده برای شناسایی تاخورده‌گی یک پروتئین مقصد از روی ساختارهای شناخته‌شده حاصل شده است. این مسئله به نام تشخیص تاخورده‌گی معروف شده است. البته این روش تنها در صورتی مؤثر است که ساختار یک یا چند پروتئین با تاخورده‌گی مشابه پروتئین مقصد در پایگاه داده شما موجود باشد.

چه کسی را باید مقاعده کرد؟ فهرست زیر به گونه‌ای مرتب شده است که تقریباً از بالا به بایین سختگیری کم می‌شود. اغلب دانشمندان پذیرفته‌اند که اقنانع «حامیان مالی طرح» از همه مهمتر است!

- | |
|----------------------------|
| چه کسی را باید مقاعده کرد؟ |
| ۱. بلورشناسان |
| ۲. متخصصان طیف‌شناسی NMR |
| ۳. حامیان مالی طرح |
| ۴. داوران مقاله‌ها |
| ۵. همکاران |
| ۶. مادرتان |

پیش‌بینی که از بانکهای داده استفاده می‌کنند عبارت‌اند از (۱) روش‌های برای هدلسازی هاستگی — پیش‌بینی ساختار مقصود از روی یک پروتئین خیلی مشابه که ساختار آن کاملاً شناخته شده است؛ و (۲) روش‌های تشخیص تاخورده‌گی — تخمین زدن تطابق‌بذری توالی آمینو اسید با مجموعه الگوهای شناخته شده تاخورده‌گی پروتئین. این روشها تا حدی (اما نه کاملاً) بدلیل رشد بانکهای داده کاملتر شده‌اند. هر چه توالیها و ساختارهای بیشتری شناخته شوند، احتمال آنکه یک پروتئین جدید مشابه پروتئین باشد که قبلاً شناخته شده است بیشتر می‌شود. بر عکس، روش‌های ابتدا به ساکن کندری پیش می‌روند. بعد از یکی از رقابت‌های اخیر CASP، یکی از اظهارنظرهای توان با دلخوری در باره این روشها این بود که حداقل «شکست از این پس تضمین شده نیست [۸]». فرد بدین ممکن است پیش‌بینی کند که رشد بانکهای داده به این معنی خواهد بود که روش‌های مبتنی بر اطلاعات راه حل‌های عملگرایانه‌ای برای چنان اکثریت بزرگی از سؤالات ارائه خواهد کرد، که اشتیاق نسبت به حمایت از توسعه روش‌های ابتدا به ساکن نقصان خواهد یافت. بنابراین مایه شرم‌ساری خواهد بود که زیست‌شناسی محاسباتی، یکی از جالب‌ترین محاسبات زیست‌شناسی را از دست داده باشد!

مراجع

- Wigner, E.P. (1960). The unreasonable effectiveness of mathematics in the natural sciences. *Communications in Pure and Applied Mathematics* 13, 1-14.
- Sankoff, D. and Kruskal, J.B., eds. (1983). *Time Warps, String Edits, and Macromolecules: The Theory and Practice of Sequence Comparison*. Addison-Wesley, Reading, Mass.
- Needleman, S.B. and Wunsch, C.D. (1970). A general method applicable to the search for similarities in the amino acid sequence of two proteins. *J. Mol. Biol.* 48, 443-453.
- Fitch, W.M. and Smith, T.F. (1983). Optimal sequence alignments. *Proc. Natl. Acad. Sci. USA* 80, 1382-1386.
- Golub, G. and van Loan, C., *Matrix Computations*. Johns Hopkins Press. Baltimore, 2nd ed. 1989.
- Young, G. and Householder, A.S. (1938). Discussion of a set of points in terms of their mutual distances. *Psychometrika* 3, 19-22. (For background and history see [7].)
- Blumenthal, L.M. (1938). *Distance Geometries. A study of the development of abstract metrics*. University of Missouri Studies 13, #2.
- New York Times*, March 25, 1997.

- Arthur M. Lesk, "The unreasonable effectiveness of mathematics in molecular biology", *Math. Intelligencer*, (2), 22 (2000) 28-37.

* آرت لسک، دانشگاه کیمیریج، انگلستان

میزان جهش‌های اتمی در محلول نشان می‌دهد که سرعت پیمایش فرساینده فضای شکلهای ممکن، چند مرتبه بزرگی کمتر از اندازه کافی خواهد بود. (این مطلب گاهی پارادوکس لویتل^۱ نامیده می‌شود). هیچ مدرکی وجود ندارد که نشان دهد مسیر تاخورده‌گی واقعاً بر وضعیت نهایی تأثیر می‌گذارد، اگرچه به لحاظ نظری این امر ممکن است. اگر حالات تاخورده‌گی دیگری ممکن بودند، اما مسیر چنان پیش می‌رفت که به یکی از آنها منجر شود، آنگاه ما پیش‌گویان مجبو بودیم رویکرد طبیعی، نه زوری، را پذیریم.

مانع پیش‌بینی ساختار کجاست؟ ما فکر می‌کنیم نیروهایی را که شکلهای طبیعی پروتئین را پایدار می‌کنند می‌شناسیم. حتی ممکن است بتوانیم تابع انرژی شکل‌گیری صریحی از مختصات بنویسیم. کاری که باید انجام دهیم کمینه کردن آن است. اما این مهم است که تشخیص دهیم که پروتئینها، به بیان ترمودینامیکی، تنها در حاشیه پایدارند. درواقع انرژی شکل‌گیری یک پروتئین تاخورده برابر با اختلاف بسیار کوچکی بین عبارات متضاد بزرگ است، چیزی که برای تحلیلگران عددی کابوس است.

آیا مشکل این است که نمی‌توانیم تابع انرژی را به اندازه کافی دقیق بنویسیم، یا تابع چنان پیچیده است که نمی‌توانیم بهینه‌اش کنیم؟ یک آزمون، کمینه کردن انرژیهای شکل‌گیری پروتئینهاست، با شروع از حالت‌های طبیعی شناخته شده آنها. چنین محاسباتی به سمت شکلهای با کمترین انرژی نزدیک به نقطه شروع همگرا می‌شوند، و نشان می‌دهند که تابع انرژی در همسایگی پاسخ صحیح رضایت‌بخش است. (این مطلب چندان شکفت‌انگیز نیست، زیرا تابع براساس تنظیم پارامترها تعریف شده‌اند به‌گونه‌ای که حالات طبیعی مشاهده شده را باز تولید کنند). اما این کافی نیست. تابعی که در همسایگی نقطه مینیمم صحیح است الزاماً مجموعه کاملی از مسیرهای پیموده شده در فضای شکلهای را — که یک برنامه را قادر می‌سازند تا با شروع از یک نقطه دلخواه مینیمم سراسری را بیابد — در اختیار ما قرار نمی‌دهد.

دو مسئله وجود دارد. نخست اینکه بسیاری از نیروهایی که پروتئینها را پایدار می‌کنند کوتاه‌برد هستند. حتی اگر تابع انرژی را دقیقاً می‌دانستیم، و اگر کمینه‌سازی را از یک شکل به طور تصادفی گسترش یافته نافرده شروع می‌کردیم، ممکن بود به این نتیجه برسیم که هیچ نیروی بلندبردی وجود ندارد که سیستم را به سمت ساختار صحیح براند. دوم اینکه، حتی اگر با فروریزی سیستم به یک حالت فشرده برسیم، نمای انرژی به عنوان تابعی از مختصات چندین مینیمم موضعی را شامل می‌شود که توسط حایلهای بلند از هم جدا شده‌اند. بسیاری از این مینیمم‌های موضعی کاندیداهایی برای حالت طبیعی خواهند بود. پروتئینهای حقیقی به وسیله ترکیبی از (۱) «بردازش موازی» سنگین که در آن همه سازه‌ها همراهان ابعاد موضعی خود را در فضای شکل می‌کاوند، و (۲) تکامل مسیرهای تاخورده‌گی که سیستم را به سمت جواب صحیح هدایت می‌کنند، بر این مسائل غلبه کرده‌اند. کامپیوترهای ما نمی‌توانند به پردازش موازی دست یابند، تابع انرژی ما نمی‌توانند مسیرهای تاخورده‌گی بلندبرد را توجیه کنند، و الگوریتمهای ما نمی‌توانند به سادگی مینیمم سراسری یک تابع پیچیده چندمتغیره غیرخطی را بیابند. (این نوسانگر همساز، وقتی لازمش داریم کجاست؟!)

دشواری حالت پیشینی، همان‌طور که متوجه شده‌ایم، به پیدایش روش‌های تجربی مبتنی بر توالیها و ساختارهای شناخته شده منجر شده است. روش‌های

1. Levinthal