

پایگاه داده قطعات پتیدی:

نرم افزار جدید برای کمک به طراحی پروتئین های نو ترکیب

فاطمه فضلعلی^{*}، سید شهریار عرب^۲

۱- فاطمه فضلعلی fatemeh.fazlali@modares.ac.ir

۲- دکتر سید شهریار عرب sh.arab@modares.ac.ir

چکیده

پروتئین ها یکی از اجزای اصلی سلول ها هستند که معمولا از چندین قطعه تقریبا مجزا تشکیل شده اند. این قطعات می توانند در طراحی و ساخت پروتئین های نو ترکیب با عملکرد های خاص استفاده شوند. طراحی پروتئین با ویژگی های مشخص یکی از زیر مجموعه های مهندسی پروتئین است. طراحی پروتئین های نو ترکیب کاربرد های اساسی در پژوهش و صنعت دارد. دستیار طراحی پروتئین یک نرم افزار بیوانفورماتیکی برای طراحی پتیدی است. استفاده از روش های بیوانفورماتیک به کاهش هزینه های ساخت و بررسی کامل تمام محیط حل مسئله کمک میکند. برنامه ارائه شده دارای یک رابط کاربری دارد که با اتصال به یک پایگاه داده امکان جستجوی برای یافتن پتیدی با شرایط مشخص مورد نظر کاربران از میان بیش از چهار میلیون قطعه پروتئینی که از حدود ۱۹۶ هزار ساختار سوم پروتئین تولید شده اند را ممکن میسازد. این نرم افزار از طریق آدرس <http://bioinf.modares.ac.ir/software/linda> برای عموم قابل استفاده است.

کلمات کلیدی: پروتئین های نو ترکیب، مهندسی پروتئین، پایگاه داده قطعات پتیدی، طراحی پروتئین، طراحی پتیدی، داده های حجیم پروتئینی، بیوانفورماتیک، بیوانفورماتیک ساختاری

۱. مقدمه

زیست شناسی مولکولی در سال های اخیر بروی شناخت بیشتر سلول متمرکز شده است. روزانه حجم وسیع و متنوعی از داده های جدید در اثر بررسی سلول ها تولید و منتشر می شود. پروتئین ها از اصلی ترین اجزای سازنده سلول ها هستند. این درشت مولکول های پیچیده بیشتر عملکرد های اصلی سلول را برعهده دارند. پروتئین ها در ویژگی های ساختاری و عملکردی سلول و در بافت اندام های بدن ها تاثیر گذارند. پروتئین ها مکانیزم های مولکولی ضروری هستند که در تمام بخش های حیات وجود دارند. اجزای سازنده پروتئین ها آمینواسید ها هستند. بیست مدل آمینواسید مختلف داریم که با ترکیب کردن آنها می توان زنجیره های پروتئینی متفاوتی تولید کرد. [1]

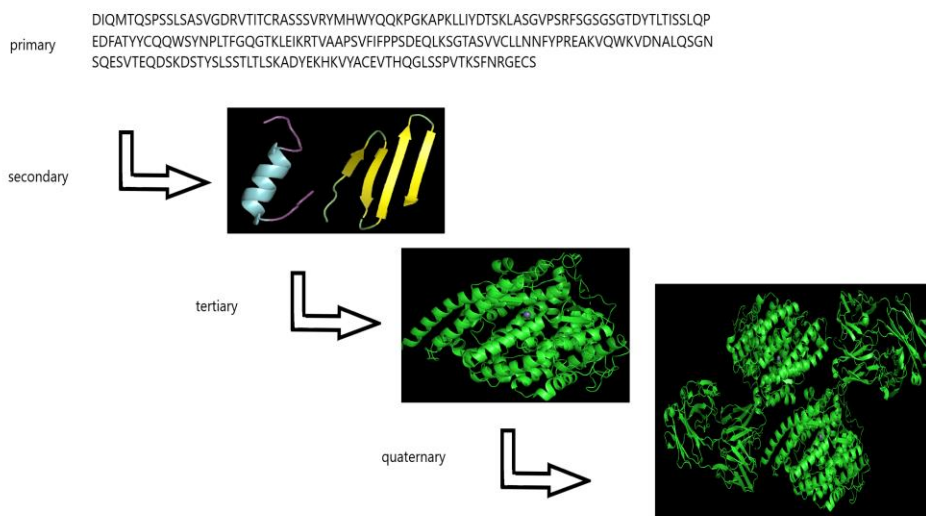
* Email: fatemeh.fazlali@modares.ac.ir

عملکرد هر پروتئین به شدت به ویژگی‌های ساختاری هر پروتئین وابسته است و با ایجاد تغییر در ساختار پروتئین‌ها می‌توان عملکردهای مختلفی در آنها ایجاد کرد. به پروتئین‌های جدیدی که از ترکیب پروتئین‌های دیگر برای ایجاد یک عملکرد خاص ساخته می‌شوند پروتئین نو ترکیب گفته می‌شود. امروزه بازار جهانی پروتئین‌های نو ترکیب بیش از یک میلیارد دلار ارزش گذاری شده است و پیش بینی می‌شود که تا سال ۲۰۳۰ به بیش از ۵ میلیارد دلار برسد. حجم و میزان رشد این بازار در یک دهه نشان از اهمیت این صنعت دارد.

مهندسی پروتئین دانشی است که به تغییر از پیش برنامه ریزی شده اسیدهای آمینه‌ها با استفاده از اطلاعات ساختاری و عملکردی پروتئین‌ها می‌پردازد. این تغییرات معمولاً بر اساس ساختار سه بعدی شناخته شده یک پروتئین معین و مکانیسم بیوشیمیایی آن است. شاخص‌ترین مزیت این کار کاهش هزینه‌های آزمایشگاهی است که با کم کردن فضای جستجو با انجام محاسبات کامپیوتری اتفاق می‌افتد. [2]

طراحی و ساخت پروتئین‌ها و آنزیم‌های جدید با عملکردهای نو بواسطه ایجاد تغییراتی در توالی‌های آمینواسیدی دارای دو راهکار عمده است. یکی طراحی هدفمند پروتئین‌ها* و دیگری تکامل جهت‌دار† در آن‌ها است. به علاوه برای ایجاد تغییرات در اندازه و ارایش فضایی مولکول‌ها هم می‌تواند با روش‌های دیگر مهندسی پروتئین انجام داد. [3]

ساختار پروتئین‡ شکل فضایی است که پروتئین به صورت خود به خودی در محیط‌های مختلف دارد. [4] پروتئین‌ها دارای چهار مدل ساختاری هستند که در شکل اول نمایش داده شده‌اند. [5]



شکل ۱- ساختارهای اول و دوم و سوم و چهارم یک پروتئین

مدل ساختاری اول توالی آمینواسیدها است که با پیوند به یکدیگر متصل شده‌اند و زنجیره پلی‌پپتید خطی را ایجاد کرده‌اند. ساختار اول ساده‌ترین و در حال مهم‌ترین ساختار یک پروتئین است زیرا ساختار فضایی و عملکردی هر پروتئین عموماً از روی ساختار اول آن تعیین می‌شود.

* Rational protein design

† Direct evolution

‡ protein structure

مدل ساختاری دوم به چیدمان فضایی آمینواسیدها در توالی‌های خطی نزدیک به هم می‌پردازد. برخی از این چیدمان‌ها قاعده مند و تکرار شونده اند. مثل مارپیچ‌های آلفا و زنجیره‌های بتا. این الگوها معمولاً با لینکرها که توالی‌های کوچک اسید آمینه هستند به یک دیگر متصل می‌شوند. با داشتن یک توالی اسید آمینه معمولاً می‌توانیم ساختار دوم این توالی را با دقت خوبی محاسبه کنیم.

مدل ساختار سوم به چیدمان فضایی آمینواسیدها در توالی‌های خطی دور از هم می‌پردازد. گاهی مشخص کردن ساختار دوم و سوم مجزا برای یک پروتئین کوچک ممکن نیست.

مدل ساختار چهارم برای پروتئین‌هایی است که بیشتر از یک زنجیره پلی‌پپتیدی دارند. در این پروتئین‌ها هر زنجیره پلی‌پپتیدی یک زیر واحد نامیده می‌شود و ساختار چهارم آرایش فضایی و نحوه تماس زیر واحد‌ها بر هم را نشان می‌دهد. در پروتئین‌های کوچک که فقط یک زیر واحد دارند ساختار سوم و چهارم یکسان است.

DSSP* عبارتی برای تعریف ساختار دوم پروتئین است. در سال ۱۹۸۳ اولین بار الگوریتمی برای ساختار دوم ارائه شد تا با استفاده از مفاهیم بصری به ساختار دوم برسیم. در واقع این الگوریتم با استفاده از مختصات اتمی سعی در شناخت الگوها دارد. هر الگو دارای ویژگی‌های ثابتی است و الگوریتم در کل هشت نوع ساختار دوم برای اسید آمینه‌ها متصور می‌شود. [6]

از سال ۱۹۷۱ پایگاه داده پروتئین‌ها[†] شروع به آرشیو اطلاعات ساختار سوم پروتئین‌های موجود کرده‌است و دسترسی به اطلاعات آن رایگان و عمومی است. این پایگاه داده هر هفته توسط مؤسسه بین‌المللی **worldwide Protein Data Bank** بروز رسانی می‌شود. [7]

هر رکورد **PDB** شامل اطلاعات شناسایی مانند عنوان، نام نویسندگان، اطلاعات آزمایش انجام شده (مثل روش و میزان تمایز)، تصویر ثابت و ساختار سه بعدی قابل دست‌ورزی پروتئینی است.

۲. روش

برای آسان‌تر کردن فرآیند طراحی پروتئین‌های نو ترکیب یک نرم افزار تحت وب طراحی شده است که به یک پایگاه داده حاوی قطعات پپتیدی متصل است. کاربران برای جستجو در میان قطعات پپتیدی می‌توانند از شش ویژگی استفاده کنند و فضای جستجو را محدود کنند. این شش ویژگی شامل تعداد و توالی آمینو اسیدهای موجود، توالی ساختار دوم هر قطعه، قطبیت هر توالی و فاصله بین ابتدا تا انتهای قطعه پپتیدی میزان دسترسی پذیری هر آمینواسیدها هستند که در ادامه توضیح داده می‌شوند. برای محاسبات این ویژگی‌ها از جدول ۱ استفاده شده است.

تعداد و توالی آمینو اسیدهای موجود: کاربران می‌توانند قطعاتی شامل ترکیبات مختلف بیست نوع آمینواسید موجود را برای یافتن لینکر مناسب به رابط گرافیکی بدهند. این عدد می‌تواند بین سه تا بیست آمینواسید باشد.

* Define Secondary Structure of Protein

† Protein Data Bank archive (PDB)

توالی ساختار دوم هر قطعه : قطعات پپتیدی می توانند هر هشت حالت مختلف در الگوریتم DSSP را برای جستجو در فضای PDB به رابط گرافیکی بدهند.

قطبیت هر توالی : توالی های آمینواسید را می توان از جنبه های مختلفی دسته بندی کرد. یکی از این دسته بندی ها چگونگی میان کنش آن ها با آب است و فهم ای عملکرد برای طراحی پروتئین های جدید ضروری است. به کمک جدول خواص آمینواسید ها (جدول یک) میزان قطبیت هر یک از قطعات پپتیدی از قبل محاسبه شده و در پایگاه داده قرار دارد و کاربران این ویژگی هم برای محدود کردن فضای جستجو می توانند استفاده کنند. هر یک از اینو اسید ها می تواند قطبی* یا غیر قطبی† باشند.

فاصله بین ابتدا تا انتهای قطعه پپتیدی: براساس واحد آنگستروم(Å)‡ می تواند به رابط گرافیکی داده شود تا تمامی قطعاتی با این طول از پایگاه داده پیدا شوند.

میزان دسترسی پذیری§ هر آمینواسید ها : برهم کنش های آبگریز نقش عمده ای در شکل گیری و پایداری پروتئین ها دارند. میزان آبدوست یا آبگریز یا آمینواسید را با میزان در دسترس حلال بودنش می سنجند. با توجه به خواص زنجیره های جانبی و قطبی بودن مولکول آب اکثر امینو اسید ها در آب حل می شوند. میزان دسترسی پذیری از قبل برای آمینواسید ها محاسبه شده و کاربران میتوانند قطعه پپتیدی مورد نظرشان را براساس این که می خواهند چه بخش هایی روی سطح پپتید(در تماس با حلال) و چه بخش هایی در داخل ساختار پپتید (بدون تماس مستقیم با حلال) باشد هم انتخاب کنند.

* Polar

† Nonpolar

‡ Angstrom

§ Accessibility

جدول ۱: جدول خصوصیات اسید آمینه های طبیعی که شامل ستون های بار الکتریکی، قطبیت، پی اچ ایزو الکتریک، شدت قدرت اسید، جرم مولکولی و نام سه و تک حرفی اسید آمینه هاست.

NO	AA_3	AA_1	Charge	Polarity	PI	PKa	PKb	PKx	MW	RW	Max_ACC
1	Ala	A	0	N	6	2.34	9.67	71.08	89.1	89.01	120.56
2	Arg	R	1	P	10.75	2.34	9.04	12.48	174.2	156.19	229.51
3	Asn	N	0	P	5.41	2.17	8.8	0	132.12	114.11	149.85
4	Asp	D	-1	P	2.77	8.8	9.6	3.65	133.11	115.09	157.04
5	Cys	C	0	P	5.07	1.88	10.28	8.18	121.16	103.15	143.79
6	Glu	E	-1	P	3.22	1.96	9.6	4.25	147.13	129.12	188.42
7	Gln	Q	0	P	5.65	2.19	9.13	0	146.15	128.13	186.83
8	Gly	G	0	P	5.97	2.17	9.6	0	75.07	57.05	89.41
9	His	H	1	P	7.59	2.34	9.17	6	155.16	137.14	200.14
10	Ile	I	0	N	6.02	1.82	9.6	0	131.18	113.16	196.42
11	Leu	L	0	N	5.98	2.36	2.36	0	131.18	113.16	206.32
12	Lys	K	1	P	9.74	2.18	2.18	10.53	146.19	128.18	213.74
13	Met	M	0	N	5.47	2.28	2.28	0	149.21	131.2	216.63
14	Phe	F	0	N	5.48	1.83	1.83	0	165.19	147.18	227.46
15	Pro	P	0	N	6.3	1.99	1.99	0	115.13	97.12	155.07
16	Ser	S	0	P	5.68	2.21	2.21	0	105.09	87.08	128.27
17	Thr	T	0	P	5.6	2.09	2.09	0	119.12	101.11	138.58
18	Trp	W	0	P	5.89	2.83	2.83	0	204.23	186.22	269.35
19	Tyr	Y	0	P	5.66	2.2	2.2	10.07	181.19	183.18	241.54
20	Val	V	0	N	5.96	2.32	2.32	0	117.15	99.13	169.82

در نهایت کاربر می تواند یکی یا چند تا از فیلد های موجود را پر کند و نرم افزار تمام قطعات پروتئینی که ویژگی های مورد نظر کاربر را دارند را پیدا و به او نمایش دهد. این قطعات می توانند در ادامه برای پیدا شدن قطعه ای که بتوانند پروتئین نو ترکیب مد نظر را بسازد مورد استفاده قرار بگیرد.

در مرحله نخست توسعه نرم افزار ابتدا نسخه ماه جولای بانک اطلاعات پروتئین ها (PDB) که دارای بیش از ۱۹۶۹۷۰ ساختار سوم پروتئین بود دریافت شد. ساختار های سوم دریافتی را فیلتر و با میزان دقت مشخص جداسازی کردیم و تعداد ۷۸۰۰۰ پروتئین را استخراج کردیم.

بسیاری از پروتئین ها به یک دیگر شبیه هستند ولی برای این نرم افزار ما ناچار به انتخاب نمونه های غیر تکراری هستیم چرا که با انتخاب پروتئین های مشابه نتایج نهایی متعلق به گروه های خاصی خواهند بود و وزن دسته ها توزیع داده ها را بر هم خواهد زد. به همین علت پروتئین هایی با تشابه ۹۵ درصد از بانک کنار گذاشته شده و تنها با ۷۸۰۰۰ پروتئین باقی مانده بانک داده پر شده است.

در مرحله ی بعد برای تبدیل فایل های اطلاعات ساختار سوم پروتئین ها (فایل های pdb) به فایل های اطلاعات ساختار دوم (فایل های dssp) یک اسکریپت پایتونی نوشته شد. (در این مرحله از اسکریپت پایتونی الگوریتم dssp استفاده شد.) اطلاعات مورد نیاز شامل نام پروتئین ها، نام زنجیره های پلی پپتیدی و محتوای اسید آمینه ها، ساختار دوم، میزان سطح در دسترس در هر پروتئین و موقعیت X,Y,Z مرکز تک تک اسید آمینه ها استخراج شد و به صورت ستون هایی در یک فایل متنی قرار گرفت .

در این مرحله فایل خروجی اسکریپت پایتونی قبلی را وارد یک پایگاه داده MySQL کردیم و به صورت یک جدول بزرگ ذخیره کردیم. سپس قطعات پپتیدی را به طول های ۳ تا ۲۰ اسید آمینه از روی جدول ساختیم. و محتوا را به صورت جدول های جدیدی که شامل ستون های توالی آمینو اسید های موجود، توالی ساختار دوم هر قطعه، قطبیت هر توالی و فاصله بین ابتدا تا انتهای قطعه پپتیدی بود قرار دادیم.

برای اتصال پایگاه داده و وب سرور اسکریپتی نوشته شد تا ارتباط بین درخواست کاربر و سیستم مدیریت پایگاه داده را ایجاد کند.

یک اسکریپت پایتونی هم برای مدیریت ورود محتوای جدید به پایگاه داده با هدف بروز نگه داشتن پایگاه داده با پایگاه اطلاعات پروتئین ها (PDB) ایجاد شد.

به این ترتیب نرم افزار پرو دآ (ProDA)* و رابط کاربری اش با هدف یافتن قطعات پپتیدی مناسب از میان بیش از 40 میلیون قطعه پپتیدی موجود در پایگاه داده ایجاد و توسعه پیدا کرد .

۳. نتیجه گیری

امروزه پروتئین ها در صنایع مختلفی از جمله بیوتکنولوژی، دارویی، بهداشتی، آرایشی، غذایی، شیمیایی و کشاورزی کاربرد دارند و تولید پروتئین های نو ترکیبی که عملکرد های جدیدی و مشخصی داشته باشند یکی از نیاز های اصلی پروسه ساخت پروتئین در این صنایع است. با در نظر داشتن این موضوع نرم افزار پرو دآ و ابزار های دیگر بیوانفورماتیکی با هدف تسهیل ساخت قطعات پروتئینی طراحی شدند تا بتواند علاوه بر کاهش هزینه های آزمایشگاهی به یافتن مجموعه ی کاملی از راه حل کمک کند.

این نرم افزار پروتئین هایی را از بانک داده پروتئین ها استخراج می کند و با پردازش آنها را به قطعات پپتیدی به طول سه تا بیست آمینواسید تبدیل می کنند و به کاربران اجازه می دهد تا با تنظیم پارامتر های ورودی رابط کاربری بین چهل میلیون قطعه پپتیدی به وجود آمده قطعه پروتئینی که شرایط مد نظرشان را دارد را پیدا کنند.

۴. مراجع

* ProDA (Protein Design Assistant)

- [1] J. E. Murray, N. Laurieri, and R. Delgoda, *Proteins*. Elsevier Inc., 2017.
- [2] P. J. Carter, "Introduction to current and future protein therapeutics: A protein engineering perspective," *Exp. Cell Res.*, vol. 317, no. 9, pp. 1261–1269, 2011, doi: 10.1016/j.yexcr.2011.02.013.
- [3] R. A. Chica, "Protein Engineering, Design and Selection," *Protein Eng. Des. Sel.*, vol. 33, p. 2020, 2020, doi: 10.1093/protein/gzaa024.
- [4] R. Zaman *et al.*, "Current strategies in extending half-lives of therapeutic proteins," *J. Control. Release*, vol. 301, no. March, pp. 176–189, 2019, doi: 10.1016/j.jconrel.2019.02.016.
- [5] C. Nelson, *Lehninger's principles of biochemistry*, vol. 53, no. 9. 2013.
- [6] Y. Zhang and C. Sagui, "Secondary structure assignment for conformationally irregular peptides: Comparison between DSSP, STRIDE and KAKSI," *J. Mol. Graph. Model.*, vol. 55, pp. 72–84, 2015, doi: 10.1016/j.jm gm.2014.10.005.
- [7] I. M. L. Saur, R. Panstruga, and P. Schulze-Lefert, "NOD-like receptor-mediated plant immunity: from structure to cell death," *Nat. Rev. Immunol.*, vol. 21, no. 5, pp. 305–318, 2021, doi: 10.1038/s41577-020-00473-z.