

برنام‌خواندگان و



تجزیه و تحلیل داده‌های RNA-seq

برای دریافت فایل‌های آموزش گام‌به‌گام، نرم‌افزارها و داده‌های مورد نیاز، به پروفایل کتاب در تارنمای انتشارات دانشگاه فردوسی مشهد به نشانی زیر مراجعه فرمایید:

press.um.ac.ir

دکتر علیرضا سیفی

عضو هیئت علمی دانشگاه فردوسی مشهد

محمد رضا رضائی

سرشناسه:	سیفی، علیرضا، ۱۳۵۶-
عنوان و نام پدیدآور:	تجزیه و تحلیل داده‌های RNA-seq / علیرضا سیفی، محمدرضا رضایی؛ ویراستار علمی امین میرشمسی کاخکی؛ ویراستار ادبی هانیه اسدیپور فعال مشهد.
مشخصات نشر:	مشهد: دانشگاه فردوسی مشهد، ۱۴۰۰.
مشخصات ظاهری:	۱۲۸ ص: مصور، جدول، نمودار.
فروست:	انتشارات دانشگاه فردوسی مشهد؛ ۷۹۰.
شابک:	ISBN: 978-964-386-483-5
وضعیت فهرست‌نویسی:	فاپا.
موضوع:	آر. ان. ا.
موضوع:	نوکلئوتیدها -- توالی
شناسه افزوده:	رضایی، محمدرضا، ۱۳۷۳-
شناسه افزوده:	میرشمسی کاخکی، امین، ۱۳۶۵- ویراستار
شناسه افزوده:	دانشگاه فردوسی مشهد، انتشارات.
رده‌بندی کنگره:	QP۶۲۳
رده‌بندی دیویی:	۵۷۲/۸۸
شماره کتابشناسی ملی:	۷۶۱۴۲۶۹

تجزیه و تحلیل داده‌های RNA-seq

پدیدآورندگان: دکتر علیرضا سیفی؛ محمدرضا رضایی
ویراستار علمی: دکتر امین میرشمسی کاخکی
ویراستار ادبی: هانیه اسدیپور فعال مشهد
مشخصات: وزیری، ۱۰۰ نسخه، چاپ دوم، زمستان ۱۴۰۴ (اول، ۱۴۰۰)
چاپ و صحافی: همیار
بها: ۲,۳۰۰,۰۰۰ ریال

حق چاپ برای انتشارات دانشگاه فردوسی مشهد محفوظ است.

مراکز پخش:

فروشگاه و نمایشگاه کتاب پردیس: مشهد، میدان آزادی، دانشگاه فردوسی مشهد، جنب سلف یاس
تلفن: ۳۸۸۰۲۶۶۶ - ۳۸۸۳۳۷۲۷ (۰۵۱)
مؤسسه کتابیران: تهران، میدان انقلاب، خیابان کارگر جنوبی، بین روانمهر و وحید نظری، بن‌بست
گشتاسب، پلاک ۸ تلفن: ۶۶۴۸۴۷۱۵ (۰۲۱)
مؤسسه دانشیران: تهران، خیابان انقلاب، خیابان منیری جاوید (اردیبهشت) نبش خیابان نظری، شماره ۱۴۲
تلفکس: ۶۶۴۰۰۲۲۰ - ۶۶۴۰۰۱۴۴ (۰۲۱)

<http://press.um.ac.ir>

Email: press@um.ac.ir



فهرست مطالب

پیشگفتار	۷
فصل ۱. تاریخچه پیشرفت‌ها در ژنتیک، بیولوژی مولکولی و بیوانفورماتیک	۹
نسل دوم روش‌های توالی‌یابی	۱۲
نسل سوم روش‌های توالی‌یابی	۱۵
ابزارهای آنالیز داده‌های NGS	۱۶
منابع	۱۹
فصل ۲. آشنایی با اصول کامپیوتر برای پژوهشگران علوم زیستی	۲۱
انتخاب سیستم کامپیوتری مناسب برای آنالیز داده‌های NGS	۲۲
حافظه‌های ذخیره‌سازی ثانویه	۲۴
ریزپردازنده مرکزی یا CPU	۲۵
حافظه اصلی یا RAM	۲۶
منابع	۲۶
فصل ۳. سیستم‌عامل لینوکس	۲۷
تاریخچه نرم‌افزارهای با دسترسی آزاد و پیدایش لینوکس	۲۷
نصب سیستم‌عامل Ubuntu	۲۸
دستورات در Ubuntu	۲۹

۳۳	کلیدهای ترکیبی در خط فرمان لینوکس
۳۴	منابع
فصل ۴. محیط نرم‌افزاری و زبان برنامه‌نویسی R	
۳۵	فصل ۴. محیط نرم‌افزاری و زبان برنامه‌نویسی R
۳۷	نصب R
۳۹	توابع R
۴۰	ساده‌سازی با ایجاد متغیر
۴۰	پوشه کاری و تغییر آن
۴۱	داده‌ها در R
۴۱	ساختار داده‌ها در R
۴۵	روش‌های وارد کردن داده‌ها به R
۴۶	عملیات ریاضی روی وکتورها
۴۷	مقایسات منطقی
۴۸	عملگرهای منطقی
۴۸	ترسیم نمودارهای آماری در R
۵۰	نمودار هیستوگرام
۵۰	بسته نرم‌افزاری ggplot2
۵۴	برخی دیگر از توابع مهم R
۵۵	نوشتن برنامه در R
۵۶	منابع
فصل ۵. RNA-seq: انقلابی در مطالعه ترنسکریپتوم	
۵۷	فصل ۵. RNA-seq: انقلابی در مطالعه ترنسکریپتوم
۵۷	روش‌های ارزیابی بیان ژن
۵۹	طراحی آزمایشات RNA-seq
۶۳	روش‌های نوین RNA-seq
۶۴	استفاده از داده‌های RNA-seq موجود در NCBI

منابع ۶۶

فصل ۶. ارزیابی کیفیت و پردازش داده‌های ایلومینا ۶۷

فایل fastq ۶۷

کنترل کیفیت خوانش‌ها با استفاده از ابزار FASTQC ۷۰

نصب نرم‌افزار FASTQC ۷۰

پردازش توالی‌ها با استفاده از نرم‌افزار Trimmomatic ۷۵

منابع ۷۹

فصل ۷. هم‌ردیف کردن خوانش‌ها با توالی مرجع ۸۱

مسیرهای مختلف آنالیز داده‌های RNA-seq ۸۱

روش‌های هم‌ردیف کردن خوانش‌ها با توالی مرجع ۸۳

هم‌ردیفی روی ژنوم مرجع با استفاده از HISAT2 ۸۴

هم‌ردیف کردن خوانش‌ها با استفاده از BWA ۸۵

جدول تعداد خوانش‌های هم‌ردیف شده ۸۷

فیلتر کردن جدول خوانش‌های هم‌ردیف شده ۸۷

منابع ۸۸

فصل ۸. آنالیز بیان ژن‌ها ۸۹

بررسی اولیه خوانش‌های هم‌ردیف شده ۸۹

شناسایی ژن‌های با بیان متفاوت با استفاده از DESeq2 ۹۱

ترسیم نقشه حرارتی ۹۲

بررسی بیان ژن‌ها با استفاده از StringTies-Ballgown ۹۵

تجزیه و تحلیل‌های پس از شناسایی ژن‌های با بیان متفاوت ۹۶

تهیه ترنسکرپتوم مرجع ۹۷

منابع ۹۸

۹۹	فصل ۹. ملاحظات آماری در تجزیه و تحلیل آزمایشات RNA-seq.....
۹۹	طرح آزمایشی در مطالعات RNA-seq.....
۱۰۰	آزمایش بدون تکرار زیستی.....
۱۰۱	نرمال سازی تعداد خوانش های هم ردیف شده.....
۱۰۳	مدل سازی خوانش های RNA-seq.....
۱۰۴	یادآوری.....
۱۰۵	یادآوری.....
۱۰۶	منابع.....
۱۰۷	پیوست: نصب لینوکس مجازی روی سیستم عامل ویندوز.....
۱۲۸	نمایه.....

پیشگفتار

فناوری‌های نوین توالی‌یابی DNA که با عنوان متداول NGS شناخته می‌شوند، تحول شگرفی در رشته‌های مختلف علوم زیستی ایجاد کرده و امکان مطالعه سریع‌تر و جامع‌تر اساس ژنتیکی و مولکولی پدیده‌های زیستی را فراهم آورده‌اند. هم‌راستا با موفقیت‌های چشمگیر سخت‌افزاری در فناوری‌های توالی‌یابی DNA، پیشرفت‌های اساسی نیز در تولید ابزارهای تجزیه و تحلیل داده‌های حاصل از این فناوری‌ها حاصل شده است.

چالش اصلی در تجزیه و تحلیل داده‌های NGS، حجم بالا و پیچیدگی این داده‌هاست که بهره‌برداری مناسب از آن‌ها نیازمند استفاده از ابزارهای بیوانفورماتیکی کارا و دقیق است. این ابزارهای بیوانفورماتیکی معمولاً توسط متخصصان علوم ریاضی، کامپیوتر و بیوانفورماتیک طراحی می‌شوند و در غالب موارد بهره‌برداری از آن‌ها برای متخصصان علوم زیستی که اطلاعات کافی از علوم کامپیوتری ندارند، چالش برانگیز است. معمولاً در مراکز تحقیقاتی پویا در دنیا یک یا چند متخصص بیوانفورماتیک عهده‌دار تجزیه و تحلیل داده‌های NGS تولیدشده توسط محققان زیست‌شناسی هستند و بنابراین تمام محققان نیازمند فراگیری این روش‌های تجزیه و تحلیل نیستند. لیکن به دلیل اینکه هنوز توسعه یافتگی مطلوبی در تشکیل تیم‌های تحقیقاتی با تخصص‌های گوناگون در مراکز تحقیقاتی ایران رخ نداده است، محققان علوم زیستی نیاز دارند که اشراف نسبی بر روش‌های تجزیه و تحلیل NGS داشته باشند.

مطالعه ترنسکرپتوم (کل محتوای RNA سلول) با استفاده از روش‌های توالی‌یابی RNA (که با عنوان RNA-seq شناخته می‌شود) یکی از قدرتمندترین روش‌های مطالعه پدیده‌های زیستی برای شناسایی مکانیسم‌های مولکولی و ژنتیکی کنترل‌کننده این پدیده‌هاست. ابزارهای بیوانفورماتیکی بسیار متنوعی برای تجزیه و تحلیل داده‌های RNA-seq ارائه شده است و ابزارهای جدید نیز با سرعت در حال توسعه و ارائه است. هدف از گردآوری کتاب حاضر فراهم آوردن مجموعه‌ای است برای محققان علوم زیستی که الزاماً اطلاعات گسترده‌ای از بیوانفورماتیک ندارند، ولی نیازمند تجزیه و تحلیل داده‌های RNA-seq هستند.

در این کتاب مباحث نظری مورد نیاز در مورد RNA-seq و اصول پایه و مقدماتی کامپیوتری توضیح داده می‌شود. سپس روش‌های استاندارد تجزیه و تحلیل داده‌های RNA-seq که در حال حاضر بیشترین کاربرد را در بین پژوهشگران دارند، معرفی می‌شوند. با استفاده از داده‌های واقعی کل مسیرهای تجزیه و تحلیل به صورت گام‌به‌گام توضیح داده می‌شوند، به نحوی که مخاطب با در اختیار داشتن این کتاب بتواند کل مسیر تجزیه و تحلیل RNA-seq را روی کامپیوتر شخصی خود تمرین کند. کلیه نرم‌افزارها، داده‌های مورد استفاده

و خروجی‌های موردانتظار از اجرای هر کدام از مراحل تجزیه و تحلیل در قالب یک لوح فشرده همراه کتاب در اختیار خوانندگان قرار می‌گیرد. تلاش شده است که مطالب کتاب و لوح فشرده همراه کتاب به نحوی تنظیم شود که علاقه‌مندان به فراگیری روش‌های تجزیه و تحلیل RNA-seq با صرف حداقل زمان، اصول این روش‌ها را دریابند و بتوانند از آن‌ها در پروژه‌های پژوهشی خود استفاده کنند.

مزیت عمده این کتاب این است که ابزارهایی که در آن معرفی می‌شوند ابزارهای استاندارد، رایگان و با متن باز هستند. استفاده از این ابزارها به پرداخت هزینه و یا عدول از اصول اخلاقی مربوط به حقوق مالکیت معنوی نیاز ندارد. نکته مهم‌تر اینکه، چنانچه تجزیه و تحلیل با این ابزارها فراگرفته شود، استفاده از ابزارهای جدیدتری که به سرعت در حال توسعه و توزیع در جامعه علمی هستند، بسیار ساده‌تر خواهد بود و پژوهشگر همواره به جدیدترین ابزارهای تجزیه و تحلیل RNA-seq دسترسی خواهد داشت.

هرچند تلاش زیادی شده است که اشکالات نگارشی و علمی در متن وجود نداشته باشد، با این حال سیاست‌گذار خواهیم بود چنانچه خوانندگان محترم اشکالات احتمالی و یا پیشنهادهای سازنده خود را برای ارتقای کیفی کتاب به نحو مقتضی به این جانب منعکس کنند.

علیرضا سیفی

زمستان ۱۳۹۹

تاریخچه پیشرفت‌ها در ژنتیک، بیولوژی مولکولی و بیوانفورماتیک

ساختار دورشته‌ای DNA در سال ۱۹۵۳ کشف و مشخص شد که مادهٔ وراثتی ماریچی است دوگانه که از جفت شدن بازهای G با C و A با T از طریق ایجاد پیوندهای هیدروژنی تشکیل شده است (۱، ۲). در سال ۱۹۷۳ اولین تلاش‌ها برای توالی‌یابی DNA به ثمر نشست و والتر گیلبرت^۱ و آلن ماکسام^۲ توانستند ۲۴ جفت باز از اپرون lac در باکتری *E. coli* را توالی‌یابی کنند (۳). در مدت چهار سال بعد روش کاراتری برای توالی‌یابی توسط فردریک سنجر^۳ ارائه شد که به نام توالی‌یابی روش سنجر^۴ شناخته می‌شود (۴). به دلیل اهمیت ابداع روش توالی‌یابی، نیمی از جایزه نوبل در شیمی در سال ۱۹۸۰ به آقای سنجر و گیلبرت اهدا شد (نیمهٔ دیگر به آقای پاول برگ^۵ برای معرفی فناوری DNA نوترکیب اهدا شد). توالی‌یابی روش سنجر نسل اول روش‌های توالی‌یابی در نظر گرفته می‌شود. توضیحات این روش در شکل ۱-۱ و انیمیشن ۱-۱ آمده است.

رشد چشمگیر پیشرفت‌ها در دنیای کامپیوتر نیز از دههٔ ۷۰ میلادی آغاز شد. تا قبل از سال ۱۹۷۰ کامپیوترها ابعاد و وزنی معادل یک یخچال خانگی داشتند و کار با آن‌ها دشوار بود. در سال ۱۹۷۷ اولین نسل از کامپیوترهای شخصی به بازار عرضه شد (Apple II، Commodore PET و TRS-80) که کار با آن‌ها برای کاربران بسیار آسان‌تر بود (۶). در سال ۱۹۷۹ با ابداع اولین نرم‌افزار آنالیز توالی‌یابی، برای اولین بار کامپیوتر در تحلیل داده‌های زیستی استفاده شد (۷).

1. Walter Gilbert
2. Allen Maxam
3. Fredrick Sanger

4. Sanger sequencing
5. Paul Berg

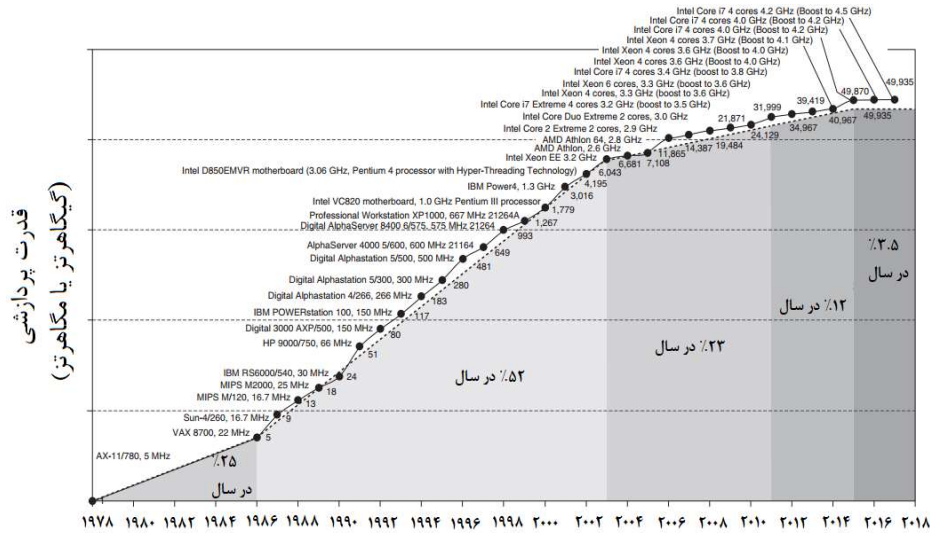
در سال ۱۹۸۴ گروه کامپیوتر و ژنتیک دانشگاه ویسکانسین^۱ مجموعه نرم‌افزاری به نام GCG را ارائه کردند که ۱۳۰ قابلیت مختلف برای آنالیز توالی‌های DNA، RNA و پروتئین داشت (۸). در همین سال‌ها بود که ریچارد استالمن^۲ بیانیه نرم‌افزار آزاد را منتشر کرد. فلسفه نرم‌افزار آزاد، منشأ چند ابتکار کلیدی در بیوانفورماتیک شد. از جمله الهام‌بخش مجموعه نرم‌افزاری متن‌باز زیست‌شناسی مولکولی اروپا شد که در ۱۹۹۶ توسعه یافته و جایگزین رایگانی برای GCG گردید. این پیشرفت در گسترش نرم‌افزارها، بدون توسعه زبان‌های برنامه‌نویسی مناسب میسر نبود. در میانه دهه ۱۹۸۰، زبان‌های برنامه‌نویسی متعددی هم‌زمان با گسترش روزافزون کامپیوترهای رومیزی، منتشر شدند که امروزه نیز هنوز در دنیای بیوانفورماتیک مطرح هستند. با ابداع زبان‌های برنامه‌نویسی جدیدتر از جمله پرل^۳ و پایتون^۴ به‌طور مؤثری پیچیدگی‌ها و مشکلات زبان‌های برنامه‌نویسی قبلی رفع شد و روند توسعه نرم‌افزارهای بیوانفورماتیک تسریع شد (۶). استفاده از کامپیوتر برای مدل‌سازی و شبیه‌سازی فرایندهای زیستی تحت عنوان زیست‌شناسی مبتنی بر کامپیوتر^۵ شناخته می‌شود که زمینه تحقیقاتی متفاوت از بیوانفورماتیک است. لازم به ذکر است که آمار زیستی^۶ نیز مبحث متفاوتی است که در آن با استفاده از روش‌ها و ابزارهای آماری، آزمایشات زیستی طراحی شده و داده‌های حاصل از این آزمایشات تجزیه و تحلیل می‌شود.

در سال ۱۹۸۸ مرکز ملی اطلاعات زیست‌فناوری^۷ (NCBI) پایه‌گذاری شد و در سال ۱۹۹۲ بانک ژن به NCBI انتقال یافت. با گسترش اینترنت در دهه ۱۹۹۰ امکان استفاده از پایگاه‌های اطلاعاتی و ثبت توالی‌های جدید در این پایگاه‌ها تسریع شد.

هرچند در دهه ۱۹۹۰ نرم‌افزارهای زیادی برای آنالیزهای بیوانفورماتیک ارائه شدند، ولی به دلیل سرعت پردازش پایین کامپیوترها در آن زمان، آنالیزهای بیوانفورماتیکی با صرف وقت زیاد قابل انجام بود. در فاصله سال‌های ۱۹۸۶ تا ۲۰۰۳ به‌طور متوسط قدرت پردازش سیستم‌های کامپیوتری ۵۲ درصد در سال ارتقا پیدا کردند (۹) (شکل ۱-۲) و در نتیجه آنالیزهای بیوانفورماتیکی تسهیل و تسریع شدند. در سال ۱۹۹۵ کرگ و نتر^۸ و همیلتون اسمیت^۹ توانستند برای اولین بار ژنوم کامل یک موجود زنده، یعنی باکتری *Haemophilus influenzae* را توالی‌یابی کنند (۱۰). یک سال بعد ژنوم یک موجود یوکاریوت تک سلولی، *Saccharomyces cerevisiae* (۱۱)، و سه سال بعد ژنوم یک موجود زنده پرسلولی، نماتد *Caenorhabditis elegans*، به‌طور کامل توالی‌یابی شد (۱۲). اولین نسخه از ژنوم انسان که توالی‌یابی آن از سال ۱۹۹۰ شروع شده بود، در سال ۲۰۰۰ ارائه شد.

1. Wisconsin
2. Richard Stallman
3. Perl
4. Python
5. Computational Biology

6. Biostatistics
7. National Center for Biotechnology Information
8. Craig Venter
9. Hamilton Smith



شکل ۲-۱ روند افزایش قدرت پردازش کامپیوترها در ۴۰ سال اخیر. قدرت پردازشی کامپیوترهای شخصی از ۵ مگاهرتز در سال ۱۹۸۷ به ۴/۵ گیگاهرتز در سال ۲۰۱۷ رسیده است. نمودار از منبع (۹) اقتباس شده است.

نسل دوم روش‌های توالی‌یابی

از نیمه دهه ۲۰۰۰ م. موج دیگری از تحولات در دنیای فناوری ایجاد شد و روش‌های توالی‌یابی نسل دوم ابداع شدند. در این روش‌ها تعداد بسیار زیادی خوانش از هر توالی تولید می‌شود. روش توالی‌یابی 454 که مبتنی بر توالی‌یابی براساس تشخیص آزاد شدن گروه پیروفسفات^۱ است، در سال ۲۰۰۵ ارائه شد. فناوری Solexa/Illumina در سال ۲۰۰۶، SOLiD در ۲۰۰۷ و فناوری Ion Torrent در سال ۲۰۱۱ به بازار فناوری معرفی شدند (۱۳). این روش‌های توالی‌یابی به‌طور کلی به‌عنوان روش‌های توالی‌یابی نسل دوم شناخته می‌شوند. در این روش‌ها قطعه DNA اولیه به روش PCR در امولسیون^۲ و یا PCR روی یک صفحه^۳ تکثیر می‌شود و سپس قطعات تولیدشده با روش توالی‌یابی براساس سنتز رشته (SBS)^۴ یا توالی‌یابی براساس اتصال (SBL)^۵ توالی‌یابی می‌شوند.

اصول بیوشیمیایی روش‌های توالی‌یابی نسل دوم با هم متفاوت است. برای مثال در روش 454 آزاد شدن پیروفسفات و در روش Ion Torrent آزاد شدن پروتون پایش می‌شود. یادآوری می‌شود که در هنگام همانندسازی DNA دو نوکلئوتید به یکدیگر متصل و پیوند فسفودی‌استر ایجاد می‌شود. در این

1. Pyrophosphate
2. Emulsion PCR
3. Bridge PCR

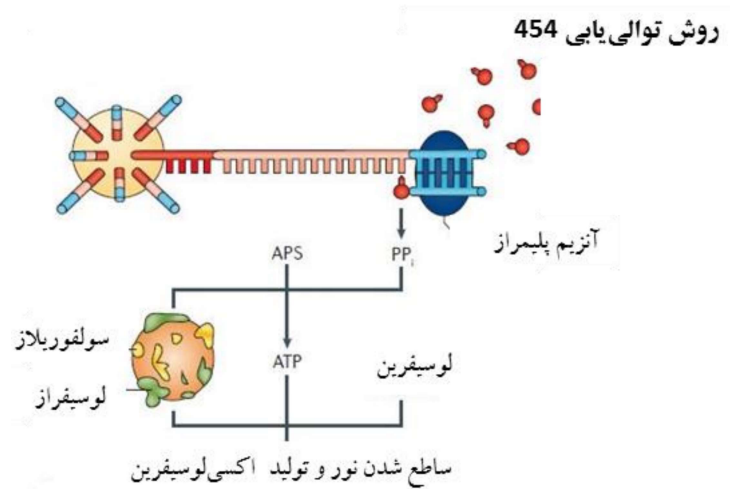
4. Sequencing by synthesis
5. Sequencing by ligation

فرایند دو مولکول فسفر (پایروفسفات، PPI) و یک پروتون نیز آزاد می‌شود. روش توالی‌یابی 454 بر تشخیص آزاد شدن PPI استوار است. روش کار به این صورت است که مخلوطی از آنزیم‌های سولفوریلاز، لوسیفراز و مادهٔ لوسیفیرین در واکنش توالی‌یابی وارد می‌شود. مولکول PPI آزادشده توسط آنزیم سولفوریلاز به ATP تبدیل می‌شود. آنزیم لوسیفراز با استفاده از این مولکول ATP لوسیفیرین را به اکسی‌لوسیفیرین تبدیل می‌کند. اکسی‌لوسیفیرین نوری ساطع می‌کند که می‌تواند توسط یک دوربین تشخیص داده شود (۱۴). روش Ion Torrent نیز مشابه روش 454 است با این تفاوت که در این روش پروتون آزادشده تشخیص داده می‌شود (شکل ۱-۳). در واقع در این روش از یک pH متر با حساسیت خیلی بالا استفاده می‌شود که می‌تواند تغییرات جزئی pH ناشی از آزاد شدن پروتون را آشکار کند.

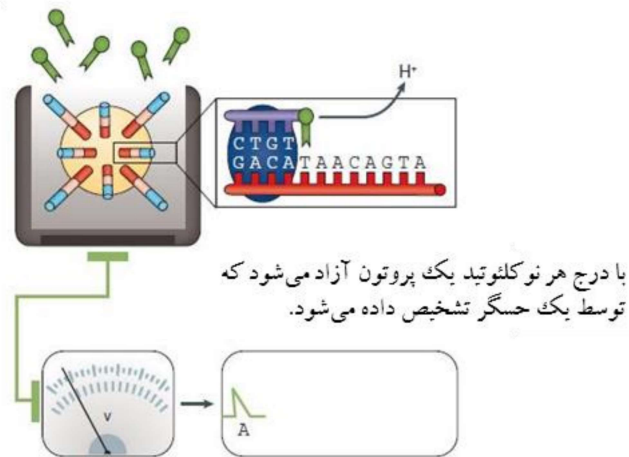
در روش توالی‌یابی ایلومینا^۱، آداپتورهای ویژه‌ای به دو انتهای مولکول‌های DNA قطعه‌قطعه شده متصل می‌شوند. سپس این مولکول‌ها به صفحه‌های ویژه‌ای که الیگونوکلئوتیدهای مکمل برای آداپتورها روی آن‌ها متصل شده است، اضافه می‌شوند. سپس یک مرحله PCR باعث ایجاد تعداد زیادی کپی از هر مولکول DNA متصل شده در هر نقطه از صفحهٔ ایلومینا می‌شود که به عنوان خوشهٔ DNA^۲ شناخته می‌شود. این فرایند تکثیر روی صفحه یا تکثیر به واسطهٔ ایجاد پل^۳ نامیده می‌شود. این نام‌گذاری به این دلیل است که در هر دور تکثیر، DNA به سمت الیگونوکلئوتید مکمل آداپتور آزاد خود روی صفحه خم و به آن متصل می‌شود. در نتیجه دو الیگونوکلئوتید متصل به صفحه به عنوان پرایمر عمل می‌کنند و تکثیر قطعهٔ جدید، روی قطعه‌ای که از یک الیگونوکلئوتید به الیگونوکلئوتید دیگر صفحه پل زده است، انجام می‌شود. سپس توالی‌یابی به روش SBS با استفاده از نوکلئوتیدهای خاتمه‌دهندهٔ برگشت‌پذیر فلوروسنت^۴ انجام می‌شود. این نوکلئوتیدها به دلیل اینکه ترکیب شیمیایی فلوروسنت آن‌ها موقعیت^۳ را اشغال می‌کند، وقتی در رشته قرار می‌گیرند از افزودن نوکلئوتیدهای دیگر جلوگیری می‌کنند، در نتیجه قبل از اینکه تکثیر ادامه یابد باید این مانع فلوروسنت برداشته شود تا توالی‌یابی ادامه یابد. در هر سیکل نور فلوروسنت ساطع شده از نوکلئوتید خاتمه‌دهنده، قبل از حذف این نوکلئوتید که برای ادامهٔ تکثیر در چرخهٔ تکثیر بعدی ضروری است، تشخیص داده می‌شود. به این ترتیب با افزودن هر یک از چهار نوکلئوتید فلوروسنت برگشت‌پذیر در هر سیکل، برخی از خوشه‌های DNA، فلوروسنت را در انتهای خود نشان می‌دهند که نشان از حضور آن نوکلئوتید در توالی آن‌هاست (۵). روش SBS در شکل ۱-۴ و انیمیشن ۱-۲ نمایش داده شده است.

1. Illumina
2. DNA cluster

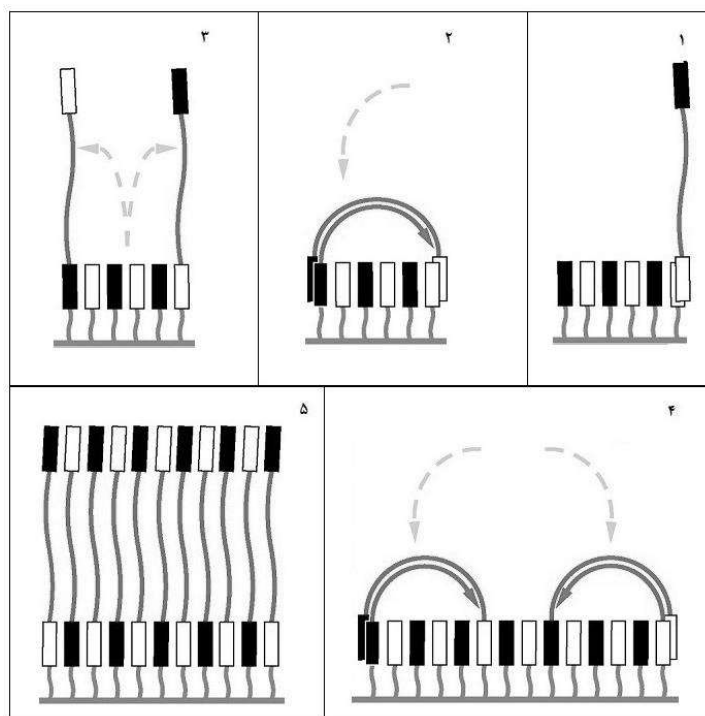
3. Bridge amplification
4. Fluorescent reversible-terminator dNTPs



روش توالی‌یابی Ion Torrent



شکل ۱-۳ مقایسهٔ اساس بیوشیمیایی دو روش توالی‌یابی نسل دوم. در نتیجه اضافه شدن یک نوکلئوتید به زنجیرهٔ درحال گسترش DNA، پروتون و گروه پایروفسفات آزاد می‌شود. در روش توالی‌یابی 454 گروه پایروفسفات و در Ion Torrent پروتون تشخیص داده می‌شود و براین اساس توالی‌یابی صورت می‌گیرد. پایروفسفات توسط آنزیم سولفوریلاز به ATP تبدیل می‌شود. مولکول ATP تولیدشده توسط آنزیم لوسیفراز استفاده می‌شود و لوسیفراز را به اکسی‌لوسیفرازین که نور قابل آشکارسازی ساطع می‌کند، تبدیل می‌کند. در Ion Torrent یک pH متر با حساسیت بالا تغییرات غلظت پروتون را آشکار می‌کند و به این ترتیب نوکلئوتید درج‌شده در توالی DNA مشخص می‌شود (۱۵).



شکل ۴-۱ روش توالی‌یابی ایلومینا. در مرحله ۱ DNA به صفحه متصل می‌شود. در مرحله ۲ DNA به سمت الیگونوکلئوتید مکمل خود خم می‌شود تا به آن متصل شود. سپس تکثیر انجام می‌شود. در مرحله ۳ دو رشته ساخته شده است. در مرحله ۴ از روی دو رشته موجود، دو رشته دیگر به روش تکثیر پل ساخته می‌شود. در مرحله ۵ خوشه‌های کلونی از قطعات DNA متصل به صفحه دیده می‌شود (۵).

نسل سوم روش‌های توالی‌یابی

در سال ۲۰۰۹ شرکت فناوری Helicos روش جدید توالی‌یابی معرفی کرد که وجه تمایز آن با روش‌های نسل دوم این بود که در آن نیاز به تکثیر DNA اولیه نبود و در تئوری امکان توالی‌یابی یک تک‌مولکول را فراهم می‌آورد. در سال ۲۰۱۱ شرکت Pacific Bioscience نیز فناوری مشابهی را روانه بازار کرد. در این روش که به روش PacBio شناخته می‌شود، مولکول DNA پلیمراز در سطح یک روزنه مستقر می‌شود و با عبور DNA از آن روزنه پلیمراز شروع به سنتز رشته مکمل آن می‌کند و به این ترتیب توالی مولکول DNA مشخص می‌شود. این روش توالی‌یابی در انیمیشن ۱-۳ نمایش داده شده است.

جدول ۱-۱ مقایسه طول خوانش و تعداد خوانش در روش‌های مختلف توالی‌یابی (۱۶)

فناوری	طول خوانش (جفت باز)	بازده (تعداد خوانش در هر اجرا)
Roche 454	۷۰۰	۷۰۰ هزار
Illumina HiSeq	۳۰۰	۳۰۰ میلیارد
SOLiD	۱۰۰	۱۰۰ میلیارد
Ion Torrent	۲۰۰	۶۰ میلیارد
PacBio RS II	۱۴۰۰۰	۴۷ هزار

فناوری جدیدتر و متفاوت‌تر توسط شرکت‌های Oxford Nanopore و IBM Transistor ابداع شد. اساس کار در این روش این است که مولکول DNA از روزنه بسیار باریکی که در آن شدت جریان الکتریکی مشخصی برقرار است، عبور داده می‌شود و تغییرات در شدت جریان که به نوع نوکلئوتید عبوری بستگی دارد، ثبت می‌شود. براساس تغییرات در شدت جریان، توالی DNA مشخص می‌شود. این روش در انیمیشن ۴-۱ نمایش داده شده است.

روش‌های توالی‌یابی نوین که سرعت و بازده بالایی دارند، انقلابی در علوم زیستی، به‌ویژه در پزشکی و کشاورزی به وجود آورده‌اند. این فناوری‌ها حجم بسیار زیادی داده ژنتیکی فراهم می‌آوردند که مطالعه سیستم‌های زنده را در سطح تک‌سلول، بافت، موجود زنده و حتی یک اکوسیستم فراهم می‌آورد. در پزشکی روش‌های تشخیصی مبتنی بر داده‌های NGS در دنیا رایج است و خوشبختانه در ایران نیز در حال متداول شدن است. در علوم کشاورزی با اتکا به این فناوری‌ها نقشه‌یابی و شناسایی ژن‌های با ارزش که کیفیت و کمیت محصولات گیاهی و دامی را افزایش می‌دهند، تسهیل و تسریع شده است. از طرف دیگر، این فناوری‌ها چالش‌های جدید و بزرگ‌تری را از نظر آنالیز داده‌های حاصل از این فناوری‌ها به همراه آورده است. قابل توجه اینکه ابزارهای بیوانفورماتیکی که در سال ۲۰۱۷ منتشر شده‌اند پنج برابر بیشتر از ابزارهای منتشرشده بین سال‌های ۱۹۹۰ تا ۲۰۰۰ بوده است (۱۷) که نشان‌دهنده لزوم یافتن روش‌ها و ابزارهای کارآمدتر برای آنالیز داده‌های NGS است.

توالی‌یابی نوین برای اهداف متنوعی از جمله توالی‌یابی کل ژنوم^۱، توالی‌یابی مجدد ژنوم^۲، توالی‌یابی ترانسکریپتوم^۳، توالی‌یابی متاژنوم^۴ و متاترنسکریپتوم^۵ استفاده می‌شوند. در همه این کاربردها مسیر کلی انجام کار تقریباً مشترک است (شکل ۱-۵).

ابزارهای آنالیز داده‌های NGS

متناسب با سرعت پیشرفت‌ها در توسعه روش‌های نوین توالی‌یابی، محققان علوم کامپیوتر، ریاضی و

1. Whole Genome Sequencing
2. Genome Re-sequencing
3. Transcriptome

4. Metagenome
5. Metatranscriptome