

# استفاده از توالی یابی طولانی برای روشن کردن مسائل اساسی و پیچیده فارماکوژنومیک

## خلاصه

استفاده از فارماکوژنومیک در عمل بالینی در حال تبدیل شدن به استاندارد مراقبت است. با این حال، به دلیل ترکیب ژنتیکی پیچیده فارماکوژن ها، در حال حاضر همه واریانت های ژنتیکی در نظر گرفته نشده است. در اینجا، ما کاربرد توالی یابی طولانی را برای حل مسائل فارماکوژن های پیچیده با آنالیز یک نمونه با مشخصه نشان می دهیم. این داده ها از خواندن های طولانی تشکیل شده است که برای حل هاپلوک های فازی پردازش شده اند. ۷۳ درصد از فارماکوژن ها به طور کامل در یک هاپلو بلوک فازی پوشانده شدند، از جمله ۱۵/۹ ژن که ۱۰۰ درصد پیچیده هستند. دقت فراخوانی متغیر در فارماکوژن ها بالا بود، با فراخوانی ۹۹.۸٪ و دقت ۱۰۰٪ برای SNV ها و ۹۸.۷٪ دقت و ۹۸.۰٪ برای ایندلز. برای اکثر تداخلات ژن-دارو در دستورالعمل های DPWG و CPIC، ژن های مرتبط می توانند به طور کامل حل شوند (به ترتیب ۶۲٪ و ۶۳٪). این یافته ها با هم نشان می دهند که داده های توالی یابی طولانی مدت فرصت های امیدوارکننده ای را برای روشن کردن فارماکوژن های پیچیده و فازبندی هاپلوتیپ ارائه می کنند و در عین حال فراخوانی دقیق را حفظ می کنند.



نیوشاده رویه<sup>۱</sup>

۱- کارشناسی ارشد ژنتیک، دانشگاه آزاد، تهران، ایران  
پژوهشگر مرکز تحقیقات پزشکی شخصی آمیتیس ژن

## مقدمه

فارماکوژنومیک (PGx) برای فردی کردن دوزهای دارو و در نتیجه بهبود نتایج درمان دارویی بسیار مهم است. PGx به فنوتیپ های استنباط شده بر اساس واریانت های شناخته شده در فارماکوژن ها متکی است. با این وجود، همه واریانت های ژنتیکی در پاسخ به دارو و فعالیت آنزیمی را نمی توان با سنجش های ژنتیکی معمول PGx، به دلیل عوامل متعدد، توضیح داد. اول، ژنوتایپینگ فعلی قادر به حل کامل ساختار ژنتیکی همه ژن های دخیل در پاسخ دارویی نیستند. دوم، مکانیسم اثر دارو و یا مسیر متابولیکی آن همیشه به طور کامل شناخته نشده

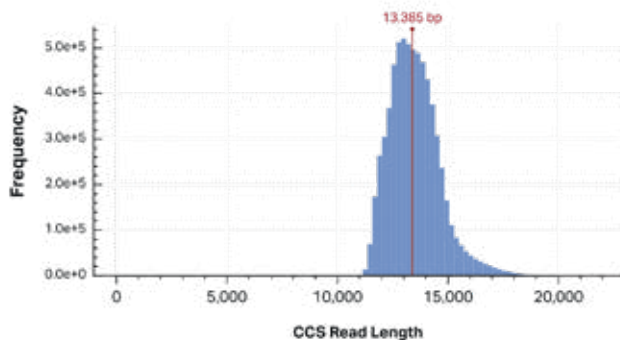
مدت در ژن FMR1 مرتبط با سندرم X شکننده و در حل ژن PKD1 برای شناسایی جهش‌های مرتبط با بیماری کلیه پلی کیستیک استفاده شده است. در نهایت، توالی‌یابی طولانی، مرحله بندی هاپلوتیپ را بدون نیاز به رویکردهای محاسباتی و یا اطلاعات خانوادگی تسهیل می‌کند. این می‌تواند در PGx اهمیت حیاتی داشته باشد که منجر به پیش‌بینی دقیق‌تر فنوتیپ می‌شود. ترکیب پیچیدگی PGx و مرحله بندی هاپلوتیپ نشان می‌دهد که توالی خوانی طولانی مدت پتانسیل بهبود قابل ملاحظه‌ای را در توانایی ما برای پیش‌بینی صحیح فنوتیپ‌های متابولیزکننده دارو دارد. در این مقاله اثبات مفهوم، ما پتانسیل توالی‌یابی طولانی مدت PacBio را برای حل نواحی پیچیده PGx با استفاده از داده‌های توالی‌یابی موجود از ژنوم به خوبی مشخص شده در نمونه مرجع بطری (GIAB) HG002 ارزیابی می‌کنیم.

## نتایج

### توضیحات داده‌ها

داده‌های توالی‌یابی قبلی منتشر شده از نمونه GIAB HG002 به خوبی مشخص شده به دست آمد. این داده شامل ۶,۷۲۸,۱۲۳ خوانش با طول میانه ۱۳.۴ kb است که ۹۷.۵ درصد از ژنوم را با میانگین پوشش نقشه برداری ۲۸ برابر پوشش می‌دهد (شکل ۱). تقریباً ۵ میلیون واریانت ژنتیکی با استفاده از HaplotypeCaller (Genome Analysis Toolkit) و DeepVariant شناسایی شد.

شکل ۱: خوانش توزیع طول. توزیع خواندن طول ژنوم در نمونه بطری HG002 پس از توالی‌یابی بر روی پلت فرم توالی Pacific Bioscience و ساخت توالی اجماع دایره‌ای.



است. برای ارزیابی اینکه چه بخشی از واریانت ژنتیکی است و چه بخشی را می‌توان توسط عوامل دیگر توضیح داد، ضروری است که بتوانیم تمام مولفه‌های ژنتیکی را که باعث پاسخ دارویی متغیر می‌شوند توضیح دهیم. با این حال، این مورد به چالش کشیده می‌شود زیرا اکثر فارماکوژن‌ها حداقل تا حدی در مناطق پیچیده ژنومی قرار دارند یا دارای انواعی مانند تکرارهای پشت سر هم و ترکیبات هیبریدی شبه ژن هستند. فن‌آوری‌های ژنوتایپینگ در حال حاضر بر اساس ریزآرایه‌های SNV (Single Nucleotide Variant) یا توالی‌یابی کوتاه خوانده می‌شوند. هر دو رویکرد در توصیف این مناطق پیچیده محدود هستند، زیرا آنها قادر به حل و فصل مناسب و قابل اعتماد مناطق بسیار همولوگ و شناسایی واریانت‌های PGx نیستند. علاوه بر این، با فزاینده هاپلوتیپ می‌توان تعیین کرد که آیا ورته‌ها روی یک آلل قرار دارند یا روی آلل‌های مختلف قرار دارند، که به طور بالقوه منجر به تفاوت‌هایی در انتساب فنوتیپ می‌شود. در حال حاضر، دیپلوتیپ‌های PGx بر اساس عدم تعادل پیوندی مرحله بندی می‌شوند. در حالی که این منجر به هاپلوتیپ‌های دقیق در مقیاس جمعیتی می‌شود، اما همیشه به فرضیات دقیق در سطح فردی منجر نمی‌شود. تأثیر این چالش‌ها در عمل بالینی زیاد است. به عنوان مثال، ژن پیچیده CYP2D6 در متابولیسم ۲۰ تا ۳۰ درصد داروهای معمول تجویز می‌شود و نمی‌توان آن را به طور کامل با توالی خوانی کوتاه مشخص کرد.

در سال‌های اخیر، فناوری‌های توالی‌یابی طولانی مدت از Oxford Nanopore و PacBio نشان داده‌اند که قادر به توصیف مناطق پیچیده ژنومی (دارویی) هستند. برای این مناطق، خواندن طولانی و با کیفیت بالا به طور قابل توجهی دقت فراخوانی انواع را بهبود می‌بخشد و امکان تفکیک دیپلوتیپ‌های کاملاً فازی را فراهم می‌کند.

ارزش توالی‌یابی طولانی مدت برای اهداف تشخیصی بیماری قبلاً نشان داده شده است. نشان داده شده است که توالی‌یابی PacBio می‌تواند CYP2D6 را با پوشش کل مکان ژن در یک خواندن طولانی با کیفیت بالا مشخص کند. اخیراً، توالی‌یابی طولانی نیز برای ژن‌های HLA در رابطه با PGx اعمال شده است. علاوه بر این، کاربرد آن در بسیاری از سنجش‌های تحقیقاتی تشخیصی بالینی چالش برانگیز مانند تکرارهای طولانی

در فارماکوژن ها نشان می دهد که با استفاده از داده های توالی خوانی طولانی در مقایسه با معیارهای فعلی، دقت از دست نمی رود، در حالی که تشخیص واریانت های پیچیده ژنتیکی را بهبود می بخشد.

### فازبندی هاپلوتیپ و هاپلوبلاک ها

با استفاده از واتس هپ، خواندن ها به صورت مرحله بندی شده و بر اساس همه واریانت های شناسایی شده به هاپلوبلاک ها تبدیل شدند. هر هاپلوبلاک یک امتداد توالی کاملاً مرحله ای را توصیف می کند که امکان توصیف کامل آن ناحیه را فراهم می کند که نشان دهنده یک آلل مادری یا پدری است. قابل ذکر است، ۷۱.۲ درصد از ژنوم را میتوان به ۱۶۱۹۳ هاپلو بلوک با طول کل هاپلو بلوک ۲.۳ میلیارد جفت باز و اندازه هاپلو بلوک میانگین ۴۰۳۰۲ جفت باز (محدوده: ۲.۹-۱ میلیون جفت باز) تبدیل کرد. یک تمایز واضح در اندازه هاپلوبلاک بین مناطق بین ژنی (میانگین ۱۴۹۶۰ جفت باز) و ویژگی های ژنکد (میانگین ۵۶۷۴۳ جفت باز)، شکل 2A مشاهده شد. اکثریت قریب به اتفاق ویژگی های Gencode به طور کامل به هاپلوبلاک ها تبدیل شدند. به طور خاص، ۷۱٪ از همه ویژگی های کدگذاری پروتئین می توانند کاملاً مرحله بندی شوند (۹۰٪ ≤) و ۲۲٪ اضافی تا حدی مرحله بندی شدند در حالی که ۷٪ حل نشده باقی ماندند (≥ ۱۰٪ فازی). الگوهای مشابهی برای سایر ویژگی های ژنکد مشاهده شد. به نظر نمی رسد طول خواندن عامل محدود کننده اصلی در حل هاپلوتیپ ها باشد زیرا درصد یک ویژگی تحت پوشش در هاپلوبلوک ها مستقل از طول ویژگی است (شکل 2B). علاوه بر این، اکثر هاپلوبلاک ها (۵۷.۷٪) از میانگین طول

### دقت و بررسی با دقت بالا در جستجو واریانت ها

برای ۱۰۰ فارماکوژن انتخاب شده، دقت و یادآوری در مقایسه با مجموعه حقیقت معیار GIAB v3.3 تعیین شد. برای SNV ها، GATK HaplotypeCaller و DeepVariant به دقت و فراخوانی مشابهی بالای ۹۹.۸ درصد دست یافتند (جدول ۱). با این حال، بررسی DeepVariant عملکرد بسیار بهتری در تشخیص ایندل (< ۹۸٪) در مقایسه با GATK (دقت: ۹۴.۵٪ و یادآوری: ۸۶.۱٪) به دست آورد. هنگام مقایسه با نتایج گسترده ژنوم گزارش شده توسط ونگر دقت و یادآوری در تشخیص انواع در فارماکوژن ها برتر است. هنگام طبقه بندی نتایج در مناطق پیچیده، دقت بالا باقی ماند، با یادآوری و دقت بیش از ۹۵٪ برای همه مناطق برای Indels و SNVs. برای بررسی GATK، دقت کمتر بود (۱۰۰-۸۵٪) در مقایسه با ۹۷-۱۰۰٪ برای بررسی DeepVariant). کاهش دقت را می توان به عملکرد پایین تر برای تکرارهای پشت سر هم و همپلیمرها نسبت داد. برای ارزیابی دقت جستجو SV در فارماکوژن ها، جستجوهای SV با معیار SV تنظیم شده برای همه SV های بیش از ۵۰ جفت باز مقایسه شد. با این حال، مناطق GIAB با اطمینان بالا همه ۱۰۰ ژن را پوشش نمی دهند. ۴۶ ژن حذف شدند، ۱۲ ژن تا حدی و ۴۲ ژن به طور کامل با داده های انتخاب شده GIAB همپوشانی داشتند. در مجموع، ۲۲ SV (بیش از ۵۰ جفت باز) در ۵۴ فارماکوژن در مقایسه با ۲۳ فهرست بندی شده در مجموعه معیار شناسایی شد. دو نتیجه منفی کاذب و یک نتیجه مثبت کاذب در نظر گرفته شد. با هم، ارزیابی عملکرد تشخیص SVs در مناطق PGx منجر به یافتن ۹۱.۳٪ و دقت ۹۴.۵٪ شد. یادآوری و دقت بالا

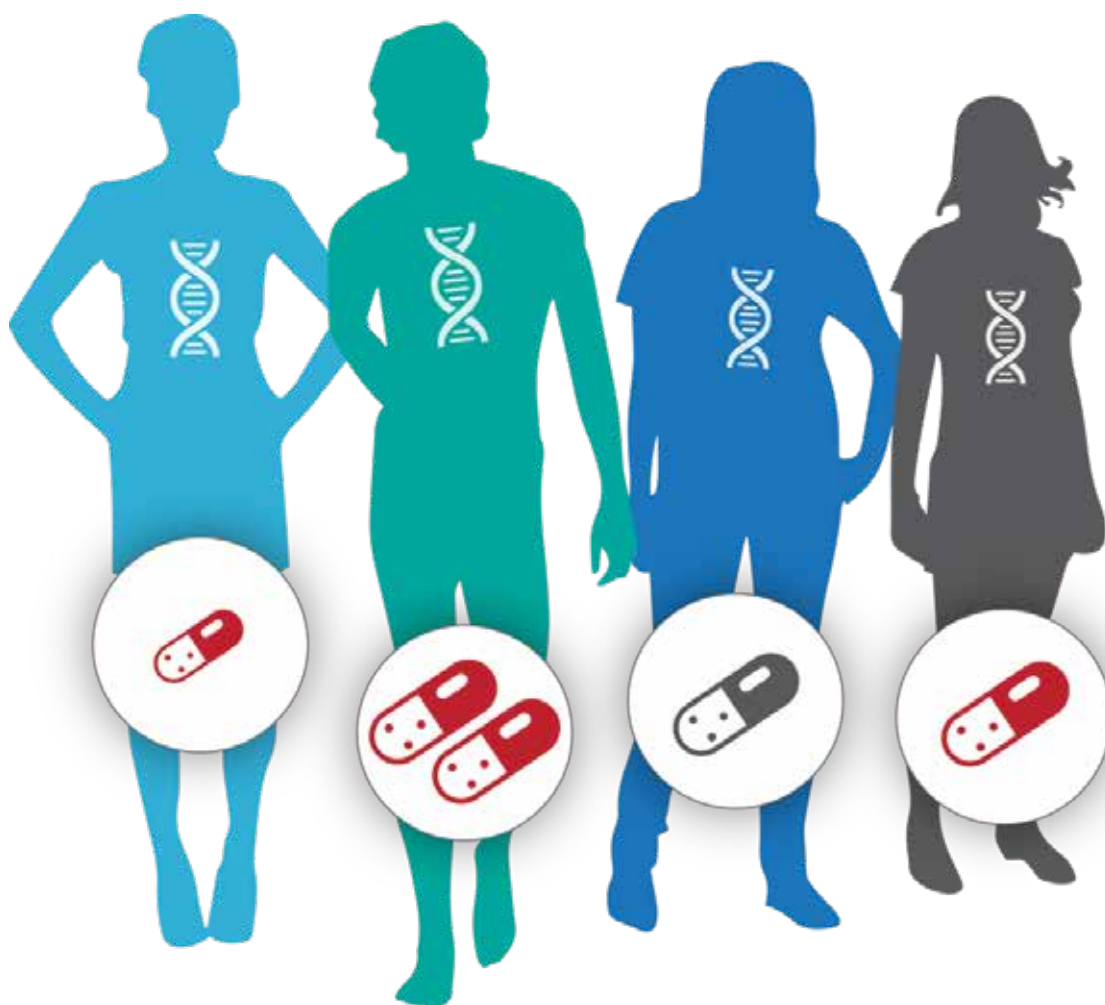
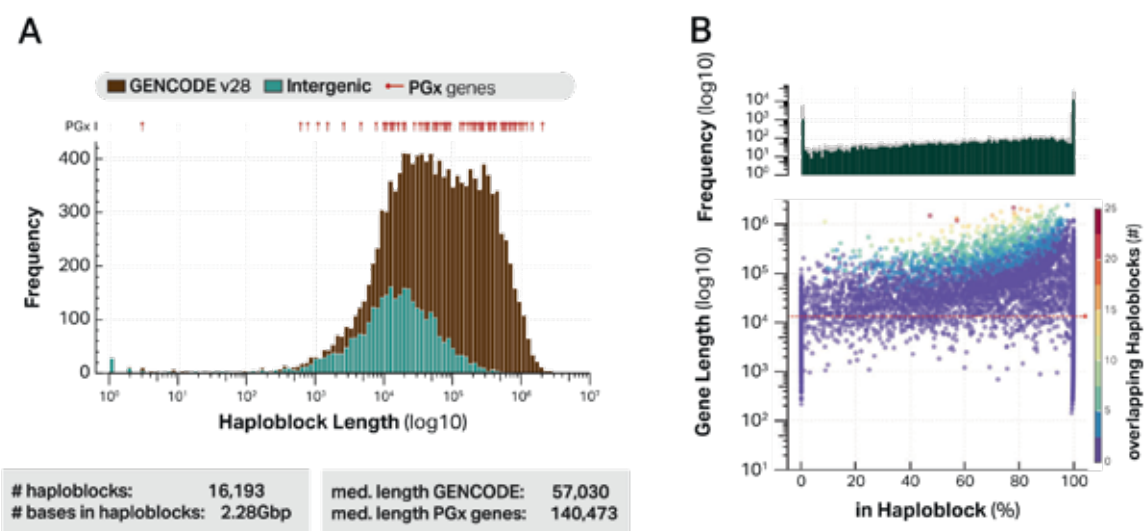
جدول ۱ عملکرد جستجوی واریانت ها برای فارماکوژن ها.

جدول ۱ عملکرد جستجوی واریانت ها برای فارماکوژن ها.					
Indels			SNVs		
F1 (%)	جستجو (%)	دقت، درستی (%)	F1 (%)	جستجو (%)	دقت، درستی (%)
۹۰.۱۰	۸۶.۱۲	۴۷.۹۴	۹۹.۹۲	۹۹.۹۶	۹۹.۸۸
۹۸.۳۷	۹۸.۰۰	۹۸.۷۴	۹۹.۹۲	۱۰۰.۰	۹۹.۸۴

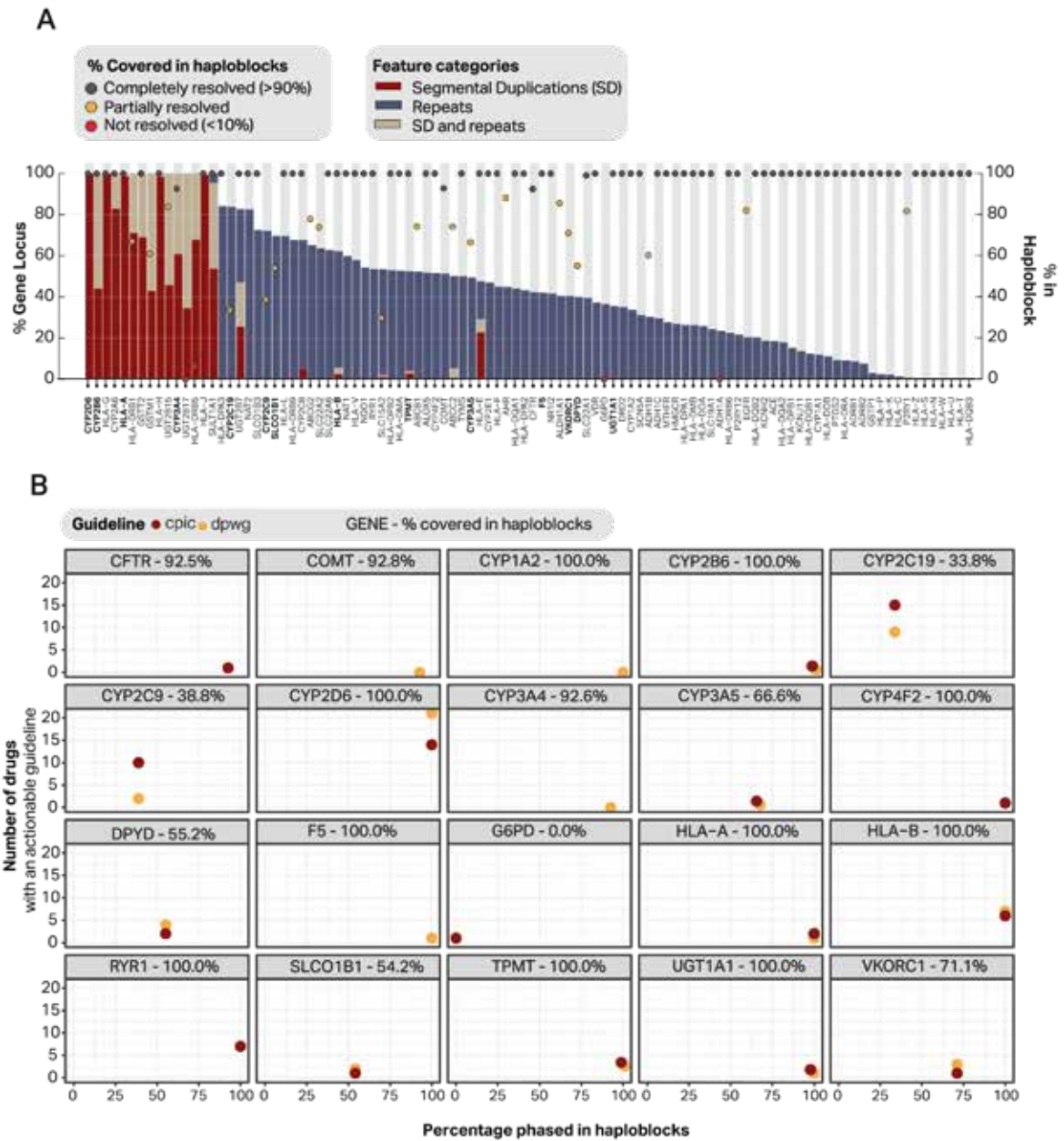
با معیار GIAB نسخه ۳.۳.۲ اندازه گیری شد. با استفاده از جستجوگر واریانت GATK و DeepVariant. واریانت تک نوکلئوتیدی SNV، درج و حذف Indels، جعبه ابزار آنالیز ژنومی GATK، توالی اجماع دایره ای CCS.

شکل ۲: وضوح هاپلوبلاک و ویژگی های GENCODE.

A: توزیع طول هاپلو بلوک طبقه بندی شده توسط ویژگی های ژنکد و مناطق بین ژنی، همپوشانی با فارماکوژن ها با رنگ قرمز برجسته شده است. B: برای هر ویژگی کدکننده پروتئین، درصدی که در مقایسه با طول ویژگی به هاپلوبلاکها تفکیک شد. خط قرمز میانگین طول خوانش را نشان می دهد. اکثر هاپلوبلاکها بزرگتر از میانگین طول خواندن هستند، که نشان می دهد نه طول خواندن، بلکه تعداد انواع هتروزیگوت برای طول یک هاپلوبلاک تعیین کننده است.



شکل ۳: پیچیدگی فارماکوژن ها و نسبت حل شده در هاپلو بلوک ها. در (A)، فارماکوژن ها و پیچیدگی آنها به درصد پوشش داده شده در هاپلوبلاک ها مربوط می شود. در ژن های پررنگ موجود در پاسپورت فارماکوژنومیک فراگیر (U-PGX) برای ژن های موجود در CPIC دستورالعمل های DPWG، تعداد دستورالعمل های عملی موجود با درصد هر ژنی که به صورت فازی به هاپلوکک تبدیل می شود، ترسیم می شود. اقدام پذیر به عنوان دستورالعملی تعریف می شود که تغییر دوز یا تغییر دارو را توصیه می کند. برای هر ژن درصد حل شده در هاپلوبلاک ها در هدر پانل گنجانده شده است. کنسرسیونم اجرای فارماکوژنتیک بالینی CPIC، گروه کاری فارماکوژنتیک هلندی DPWG.



آن هموزیگوت هستند که منجر به بلوک‌های فازی تکه تکه می‌شود. با این حال، از آنجایی که همه مناطق توالی‌یابی شده‌اند، هنوز امکان اختصاص هاپلوتیپ‌ها و فنوتیپ‌ها با استفاده از دستورالعمل‌ها و فرضیات مرحله بندی گروه کاری فعلی فارماکوژنتیک هلندی (DPWG) و کنسرسیونم اجرای فارماکوژنتیک بالینی (CPIC) وجود دارد.

برای ارزیابی کاربرد بالینی، دیپلوتیپ‌ها و فنوتیپ‌ها بر اساس پانل واریانت از کنسرسیونم فارماکوژنومیکس Ubiquitous (U-PGx) و یک خط لوله که قبلاً توسعه یافته بود، اختصاص داده شدند. در مجموع ۱۴۱۸ واریانت در ۱۰ فارماکوژن کلیدی موجود در پانل شناسایی شد که از این میان ۳۸ وریته در پانل فنوتیپ در نظر گرفته شد. واریانت‌های مرتبط بالینی در ژن‌های CYP3A5، CYP2D6 و VKORC1 شناسایی شدند. برای CYP3A5، وریته‌ی (g.99672916 C > T) rs776746 روی هر دو آلل یافت شد که منجر به یک فنوتیپ CYP3A5\*3/\*3 و یک فنوتیپ متابولیزر ضعیف شد. برای CYP3A5 وضعیت PM غیر قابل عمل تلقی می‌شود زیرا این فنوتیپ رایج‌ترین فنوتیپ در قفقازی‌ها است. برای CYP2D6 و VKORC1 فنوتیپ استنباط شده از وریته‌ی وحشی متفاوت بود. در لوکوس CYP2D6، هر دو وریته (g.42128945 C > T) rs3892097 و (g.42130692 G > A) rs1065852 هتروزیگوت بودند. با فزبنندی، مشخص شد که واریانت‌ها روی همان آلل قرار دارند که منجر به دیپلوتیپ CYP2D6\*1/\*4 و فنوتیپ متابولایزر میانی (IM) CYP2D6 است. علاوه بر این، با توجه به وجود شبه CYP2D7 غیر عملکردی که بیش از ۹۰٪ از توالی آن را با CYP2D6 به اشتراک می‌گذارد، حذف هرگونه تداخل خوانش CYP2D7 برای تعیین دقیق هاپلوتیپ‌های CYP2D6 بسیار مهم است. قرائت‌ها به اندازه کافی طولانی بودند تا تمایز واضحی بین CYP2D6 و CYP2D7 بدون هیچ گونه قرائت‌های نقشه برداری مبهم وجود داشته باشد. همین امر برای CYP2B6 و شبه CYP2B7P آن و برای جایگاه CYP3A که همه ژن‌ها همسانی توالی بالا را دارند مشاهده شد. برای VKORC1، یک واریانت هموزیگوت (NC\_000016.10: g.31093557 G > A) شناسایی شد که منجر به فنوتیپ 1173TT و منجر به کاهش فعالیت شد. به طور کلی، این

خواندن فراتر می‌روند، که نشان می‌دهد نه طول خواندن، بلکه تعداد گونه‌های هتروزیگوت و تعداد خوانده‌های هم تراز با یک منطقه ژنومی معین، عوامل محدود کننده در ساخت هاپلوبلاک هستند.

## فارماکوژن‌ها

برای هر یک از ۱۰۰ فارماکوژن انتخاب شده، بخشی از ژن‌های واقع در یک منطقه پیچیده تعیین شد که با مجموعه‌ای که به عنوان مناطق ژنومی تعریف می‌شود و با تکرارهای قطعه‌ای (SD) یا تکرار همپوشانی دارند. در مجموع، ۱۵ فارماکوژن به عنوان ۱۰۰٪ پیچیده طبقه بندی شدند در حالی که هشت فارماکوژن هیچ همپوشانی با SDs یا تکرار نشان ندادند (شکل ۳A).

برای هر یک از ۱۰۰ لوکوس، تقریباً همه واریانت‌ها را می‌توان با دقت جستجو کرد (دقت و یادآوری > ۹۹.۸٪). فزبنندی بعدی منجر به هاپلوبلاک‌هایی با طول متوسط ۱۴۰۴۷۳ جفت باز شد که در نتیجه اکثریت (۷۳/۱۰۰) ویژگی‌ها به طور کامل به هاپلوبلاک تبدیل شدند (شکل ۳A). مهم‌تر از همه، از ۱۵ فارماکوژن طبقه‌بندی شده به‌عنوان کاملاً پیچیده، ۹ فارماکوژن می‌توانند به‌طور کامل فزبنندی شوند، ۴ تا برای حداقل ۶۰ درصد و دو فارماکوژن آخر نمی‌توانند فزبنندی شوند. از ژن‌های پیچیده ۳۵، HLA مورد از ۳۷ ژن کاملاً برطرف شدند، دو مورد باقی مانده (HLA-DRB1 و HLA-DRB5) به ترتیب ۶۰.۴ و ۶۷.۱ درصد برطرف شدند.

با این وجود، چندین فارماکوژن مهم تنها می‌توانند تا حدی به هاپلوبلاک تبدیل شوند. به عنوان مثال، G6PD، CYP2C19 و DPYD به ترتیب ۵۵ و ۳۴ درصد حل شدند. از آنجایی که G6PD بر روی کروموزوم X قرار دارد و توالی‌یابی شده مذکر است، نمی‌توان مکان را به دو آلل تبدیل کرد و در نتیجه ۰٪ از لوکوس در هاپلوبوک‌های فازی پوشانده می‌شود. برای DPYD علت در ترکیبی از طول ژن طولانی (۹۰۰۰۰۰~ جفت باز) و تعداد کم واریانتی است که منجر به کشش‌های بزرگ بدون واریانت‌های هتروزیگوت می‌شود که منجر به هاپلوبوک‌های شکسته می‌شود. برای CYP2C19، بخش بزرگی در مرکز ژن وجود دارد که برای همه واریانت‌ها هموزیگوت است. به طور خاص، در کل لوکوس CYP2C19 حدود ۵۲ واریانت نوع دارد که ۳۳ وریته‌ی

نتایج نشان می‌دهد که، طبق دستورالعمل‌های اجماع عمومی در دسترس، این فرد نیاز به تنظیم دوز برای داروهایی دارد که سوبسترای CYP2D6 و VKORC1 هستند.

### ارتباط بالینی

در مجموع، ۱۵ ژن موجود در این مطالعه در دستورالعمل‌های CPIC و یا DPWG نشان داده شده‌اند که منجر به ۵۶ و ۶۷ تعامل ژن-دارو به ترتیب برای دستورالعمل‌های DPWG و CPIC می‌شود (شکل ۳). از این تعداد ۱۰ ژن (۷/۶۶ درصد) به طور کامل در هاپلو بلوک‌های فازی حل شدند. ژن‌هایی که به طور کامل حل شدند در ۳۵ تا از برهمکنش‌های ژن-دارو در DPWG و برای ۳۵ برهمکنش ژن-دارو در CPIC نقش دارند. برای ژن‌های باقی‌مانده، هنوز هم می‌توان واریانت‌ها را به دقت شناسایی کرد و با توجه به عملکرد بالینی فعلی که از داده‌های ژنتیکی غیر مرحله‌ای استفاده می‌کند، تخصیص هاپلو تیپ را ممکن می‌سازد.

### بحث

در این مطالعه اثبات مفهوم، نشان داده‌ایم که توالی خوانی طولانی مدت جستجوهای با کیفیت بالایی را در تمام فارماکوژن‌های انتخابی به دست می‌دهد. در مقایسه با آنالیز گسترده ژنوم، نتایج برای ژن‌های PGx با توجه به دقت فراخوانی و وضوح هاپلوک‌های فازی بزرگ‌تر برتر است. علاوه بر این، اکثر فارماکوژن‌های انتخاب شده می‌توانند به طور کامل در هاپلو بلوک‌های فازی حل شوند.

بر اساس فراخوانی نوع به تنهایی، داده‌های کل ژنوم طولانی خوانده شده را می‌توان برای PGx معمولی مشابه روش استفاده از NGS استفاده کرد. علاوه بر این، توالی‌یابی طولانی مدت مزیت حل آلل‌های پدری و مادری را ارائه می‌دهد. با توجه به ماهیت چندشکلی فارماکوژن‌ها، احتمال اینکه یک فرد دارای چندین گونه در یک فارماکوژن باشد بسیار زیاد است و اهمیت فازبندی هاپلو تیپ را افزایش می‌دهد. علاوه بر این، این فراوانی زیاد انواع منجر به هاپلوبلاک‌های قابل توجهی بزرگ‌تر برای فارماکوژن‌ها در مقایسه با ویژگی‌های Gencode شد.

توالی‌یابی طولانی از نظر تشخیص SNV با توالی خوانی کوتاه قابل مقایسه است و از نظر فازبندی هاپلو تیپ و SV‌های پیچیده بهتر عمل می‌کند. فازبندی هاپلو تیپ به طور بالقوه می‌تواند تفاوت بین فنوتیپ متابولایزر میانی استنباط شده (دو نوع کوتاه کننده در یک آلل) و فنوتیپ متابولایزر ضعیف (دو نوع کوتاه کننده در آلل‌های مختلف) ایجاد کند. استراتژی‌های هاپلوتا پیننگ فعلی PGx از مرحله‌بندی محاسباتی استفاده می‌کنند که منجر به مرحله‌بندی دقیق در مقیاس جمعیتی، اما نه همیشه در سطح فردی می‌شود. از آنجایی که تنظیمات دارو در سطح فردی انجام می‌شود، دقت در مورد مرحله بندی برای یک فرد بسیار مهم است. در اینجا نشان داده‌ایم که توالی‌یابی طولانی مدت اکثر فارماکوژن‌ها را قادر می‌سازد تا بدون نیاز به داده‌های شجره‌نامه یا فازبندی محاسباتی، به طور کامل در هاپلوبوک‌ها قرار بگیرند.

توالی‌یابی طولانی همچنین توصیف کاملی از هر واریانت در لوکوس‌های انتخابی PGx، از جمله واریانت‌های ساختاری و نادر ارائه می‌دهد، همانطور که با دقت و یادآوری بالا برای SNV‌ها، Indels و SV‌ها مشخص می‌شود. به عنوان مثال، میانگین طول خوانده شده (۱۳.۴kbp) تقریباً سه برابر بزرگتر از اندازه منبع (4.4kbp) CYP2D6 است، که امکان توصیف کامل لوکوس و CNV‌های بالقوه را فراهم می‌کند. تفاوت بزرگ بین DeepVariant و GATK برای Indels را می‌توان با استفاده از داده‌های طولانی خوانده شده PacBio CCS برای آموزش جستجوگر DeepVariant توضیح داد. GATK با حالت خطای توالی خواندن کوتاه به عنوان پایه، با ۱۰۰ برابر جایگزینی بیشتر از ایندل طراحی شد. از طرف دیگر DeepVariant حالت خطا را از داده‌های آموزشی PacBio HiFi یاد گرفته است که نسبت ایندل ۳۰ برابر بیشتر از جایگزینی دارد. به طور خاص، شناسایی تکرارهای Indels و پشت سر هم با استفاده از خوانش طولانی و DeepVariant به طور قابل توجهی بهبود یافته است. این تفاوت مزیت افزوده خواندن طولانی نسبت به توالی‌یابی کوتاه را برای شناسایی واریانت‌های پیچیده بار دیگر برجسته می‌کند.

برای فرد مورد مطالعه، ۱۴۱۸ SNV در لوکوس‌های PGx بالینی انتخاب شده (۱۰ ژن) شناسایی شد که ۹۴

قادر به شناسایی دقیق واریانت‌ها هستند، اما در توانایی خود برای حل همه پیچیدگی‌ها و در رابطه با فازبندی هاپلوتیپ محدود هستند.

با این وجود، لازم به ذکر است که تمام فارماکوژن‌ها نمی‌توانند به طور کامل برطرف شوند. دلیل اصلی این امر فقدان واریانت هتروزیگوت بود تا امکان ساخت هاپلوبلوک را فراهم کند. این به نوبه خود منجر به شکسته شدن هاپلوبوک‌ها و فارماکوژن‌ها می‌شود که به طور کامل قابل حل نیستند. برای فردی که مطالعه کردیم، این اثر به ویژه برای CYP2C19 و DPYD آشکار بود. با این حال، شناسایی وریده هنوز در کل لوکوس ژن امکان پذیر بود که امکان تعیین هاپلوتیپ غیر فازی را فراهم می‌کرد. برای این ژن‌هایی که نمی‌توانند به طور کامل حل شوند، روش‌های هاپلوتیپ مرسوم مبتنی بر داده‌های توالی‌یابی غیر فازی همچنان می‌توانند به کار روند که منجر به پیش‌بینی هاپلوتیپ و فنوتیپ مطابق با عملکرد بالینی فعلی می‌شود. علاوه بر این، برای DPYD، سه مورد از چهار واریانت مرتبط با بالینی هنوز مرحله‌بندی شده‌اند، که دو مورد از آن‌ها در هاپلوبوک یکسان هستند. نشان‌دهنده عدم فازبندی کامل به این معنی نیست که هیچ یک از واریانت‌های مربوطه نمی‌تواند مرحله‌بندی شود. از آنجایی که پوشش در تمام فارماکوژن‌ها کافی بود، این عدم فازبندی ناشی از ساختار ژنتیکی افراد است، به دلیل فقدان واریانت‌های هتروزیگوت در این ناحیه، و نه با توالی‌یابی به خودی خود، این به راحتی قابل حل نیست. برای فرد دیگری، همین مشکل هاپلوبوک‌های شکسته ممکن است در ژن‌های دیگر بسته به ژنتیک آنها مشاهده شود. در حالی که توالی‌یابی طولانی برای فارماکوژنومیک بالینی امیدوارکننده به نظر می‌رسد، هزینه‌ها و زمان چرخش مرتبط با آن در حال حاضر برای تشخیص PGx با توان بالقوه بالا بسیار زیاد است. در حال حاضر، این باعث می‌شود توالی‌یابی طولانی با آرایه‌های SNV سریع مورد استفاده در PGx بالینی سازگار نباشد. با این حال، هزینه‌های توالی‌یابی به سرعت در حال کاهش است. علاوه بر این، ژنوتیپ پیشگیرانه محبوب‌تر می‌شود که باعث می‌شود زمان چرخش طولانی‌تر دیگر مسئله‌ای نباشد.

در این مطالعه از داده‌های ژنتیکی یک نمونه DNA با کیفیت بالا استفاده شد. در عمل بالینی، کیفیت

درصد آن‌ها کاملاً فازبندی شده بودند، که نشان‌دهنده فراوانی واریانت‌ها در فارماکوژن‌ها است. علاوه بر این، ماهیت مرحله‌ای این داده‌ها می‌تواند به بهبود درک ما از هاپلوتیپ‌ها و ترکیب‌های مختلف کمک کند. بنابراین، فناوری‌های توالی‌یابی طولانی مدت پتانسیل تغییر دانش ما در مورد عوامل ژنتیکی را دارند که در پاسخ دارویی متغیر نقش دارند.

قبل از اجرای توالی‌یابی طولانی در عمل بالینی، ابزارهایی برای کمک به تفسیر مورد نیاز است. چندین گروه برای توسعه چنین ابزارهای ترجمه‌ای برای PGx تلاش کرده‌اند. با این حال، هنوز محدودیت‌هایی برای این ابزار وجود دارد. اول، آن‌ها اغلب طیف وسیعی از واریانت‌های شناخته شده و هاپلوتیپ‌های مرتبط با آن‌ها را پوشش می‌دهند. با این حال، برای هر \*-هاپلوتیپ تأثیر بالینی مشخص نیست، بنابراین گاهی اوقات منجر به هاپلوتیپی می‌شود که اثر آن ناشناخته است و اجرای آن در عمل بالینی را دشوار می‌کند. ثانياً، ابزارها همیشه نتیجه یکسانی را برای یک فرد ارائه نمی‌دهند، که نشان می‌دهد مفروضاتی که این ابزارها بر اساس آن‌ها هستند قابل مقایسه نیستند. برای اینکه فقط \*-هاپلوتیپ‌های مرتبط با بالینی را در آنالیز خود لحاظ کنیم، تحلیل خود را از ابزار بالینی به پانل واریانت‌های تعریف شده توسط کنسرسیون U-PGx محدود کرده‌ایم. البته لازم به ذکر است که این امر منجر به حذف اکثر واریانت‌ها در همه لوکوس‌های PGx می‌شود، به دلیل این واقعیت که هنوز دانش کافی در مورد عملکرد این واریانت‌ها وجود ندارد. برای نشان دادن تأثیر خواندن طولانی بر PGx بالینی، ما نتایج توالی‌یابی را در چارچوب دستورالعمل‌های DPWG و CPIC ارزیابی کرده‌ایم. بر اساس واریانت‌های ژنتیکی مشاهده شده در فرد مورد مطالعه، دستورالعمل‌ها تنظیم دارو یا دوز را برای ۲۲ دارو توصیه می‌کند. از تمام تداخلات ژن-دارو در دستورالعمل‌ها (۵۳ برای DPWG و ۵۴ برای CPIC)، اکثریت قریب به اتفاق (۳۵ برای هر دو) با یک ژن پیچیده (جزئی) همراه بود که می‌توانست به طور کامل در یک هاپلوبلاک حل شود. همانطور که در این مطالعه نشان دادیم، توالی‌یابی طولانی قادر به حل این پیچیدگی‌ها و ساختن هاپلوبلاک‌های بزرگ است که امکان فراخوانی هاپلوتیپ دقیق‌تری را فراهم می‌کند. پانل‌های SNV و توالی‌یابی کوتاه، از سوی دیگر،



طولانی برای PGx کافی است. بر اساس این داده‌ها در مورد دقت فراخوانی و توانایی تفکیک فارماکوژن پیچیده به هاپلوبوک‌های فازی، نتیجه می‌گیریم که داده‌های توالی‌یابی طولانی مدت فرصت‌های خوبی برای روشن کردن لوکوس‌های پیچیده PGx و فازبندی هاپلوتیپ ارائه می‌دهد در حالی که فراخوانی دقیق وریته در فارماکوژن‌های انتخابی حفظ می‌شود.

## مواد و روش‌ها

### توضیحات داده‌ها

نتایج توالی‌یابی طولانی در دسترس عموم از نمونه GIAB HG002 با توالی‌یابی PacBio تعیین توالی شد و با استفاده از CCS (Circular Consensus Sequencing) خوانده شد. یک نمونه GIAB انتخاب شد زیرا اینها با نتایج معیار موجود بسیار خوب مشخص می‌شوند. خواندن‌های CCS با استفاده از نرم افزار CCS نسخه 3.0.0 ایجاد شد. قرائت‌های HiFi به دست آمده با استفاده از NGMLR aligner v0.2.7 با ژنوم مرجع GRCh38 تراز شد. انواع ژنتیکی با استفاده از GATK HaplotypeCaller (v.4.0.6.0) و DeepVariant (v.0.7.1) شناسایی شدند. مجموعه‌ای از ۶۴ فارماکوژن که قبلاً توسط Lauschke به همراه ژن‌های پیچیده HLA توصیف شده بود برای آنالیز PGx انتخاب شدند.

### ساخت هاپلوبلاک

واریانت‌های نامیده شده توسط GATK با استفاده از WhatsHap برای به دست آوردن واریانت‌های SNV فازی و Indel مرحله بندی شدند. از قرائت‌های مرحله‌ای، هاپلوبلاک‌ها ساخته و در فایل‌های GTF و BED ذخیره شدند. هر هاپلوبلاک با تطبیق قرائت‌های مرحله‌ای بر اساس واریانت‌های موجود در آن‌ها ساخته شد تا طول توالی قابل حل را افزایش دهد. یک هاپلوبلاک نشان دهنده یک امتداد توالی ناگسستنی است که بر اساس خوانش‌های مرحله‌ای همپوشانی دارند و زمانی متوقف می‌شود که ناحیه‌ای در ژنوم فقط با خوانش‌های بدون هیچ گونه‌ای پوشانده شود، دیگر تفاوتی در واریانت‌های بین دو آلل وجود ندارد یا اگر منطقه فاقد پوشش باشد. پس از آن، همه جایگاه‌ها به یکی از سه ویژگی طبقه‌بندی شدند: ویژگی‌های ژن کد (v28)، ژن‌های

بالا ممکن است همیشه تضمین نشود. با این وجود، کاربردهای قبلی توالی‌یابی طولانی مدت در یک محیط بالینی یا با استفاده از DNA به دست آمده بالینی منجر به نتایج با کیفیت خوبی شده است. علاوه بر این، از سال ۲۰۲۰ یک گردش کار ورودی DNA بسیار کم PacBio که تنها به ۵ اینچ از DNA نیاز دارد، در دسترس بوده است. بنابراین انتظار می‌رود که نتایج توالی‌یابی با کیفیت بالا را بتوان با نمونه‌های بالینی جمع‌آوری‌شده به‌طور معمول به‌دست آورد.

دقت و ارزش توالی‌یابی طولانی قبلاً در کل داده‌های ژنوم بررسی شده است، که ممکن است یک رویکرد هدفمند همانطور که در اینجا ارائه کردیم غیر ضروری به نظر برسد. با این حال، به خوبی ثابت شده است که پیچیدگی مناطق فارماکوژنومیک ژنوم، سنجش‌های کنونی را در حل ساختار ژنتیکی آن‌ها به خطر می‌اندازد و در نتیجه قابلیت اطمینان و کامل بودن سنجش‌های فنوتیپی را محدود می‌کند. تفاوت در ساختار ژنتیکی فارماکوژن‌ها در مقایسه با ژن‌های کد کننده پروتئین عمومی، برون یابی مستقیم از نتایج کل ژنوم را غیر قابل اعتماد می‌کند. مهمتر از همه، آنها حاوی واریانت‌های بیشتری هستند که با هم بر پاسخ دارویی تأثیر می‌گذارند. این تعداد زیاد پلی‌مورفیسم‌ها منجر به این فرضیه می‌شود که فارماکوژن‌ها به دلیل فراوانی بیشتر واریانت‌های هتروزایگوت، آسان‌تر فازبندی می‌شوند، همانطور که در مطالعه ما تأیید شد. در واقع، دقت در فارماکوژن‌ها بیشتر از سایر ژن‌ها بود، در حالی که خواندن کوتاه دقت بسیار پایین‌تری در تشخیص واریانت‌های ژنتیکی در این مناطق پیچیده دارد. توانایی توالی‌یابی طولانی مدت برای حل فارماکوژن‌ها قبلاً در مطالعات توالی‌یابی هدفمند نشان داده شده بود. با این حال، این مطالعه با هدف ارائه یک نمای کلی جامع از کاربرد توالی‌یابی طولانی مدت در حل فارماکوژن‌های پیچیده و اطلاع رسانی در مورد مناطقی بود که همچنان چالش برانگیز هستند.

این مطالعه به داده‌های با کیفیت بالا از یک موضوع محدود می‌شود و به عنوان اثبات مفهومی برای کاربرد توالی‌یابی طولانی در PGx عمل می‌کند. علیرغم این محدودیت، ما احساس می‌کنیم که این به عنوان یک مطالعه اثبات مفهوم برای بررسی پتانسیل توالی‌یابی

همپوشانی دارند به عنوان «پیچیده» تعریف می‌شود.

### ارتباط بالینی

یک پایپ‌لاین که قبلاً توسعه یافته بود برای اختصاص هاپلوتیپ‌ها و فنوتیپ‌ها به فارماکوژن‌های مرتبط بالینی بر اساس دستورالعمل‌های DPWG استفاده شد. ژن‌های انتخاب شده بر اساس پانل کنسرسیون U-PGx و شامل ۱۰ فارماکوژن کلیدی و ۳۸ وریته بودند. این پایپ‌لاین که در اصل برای داده‌های NGS طراحی شده بود، شامل ژن‌های UGT1A1 و HLA-B که در پنل کنسرسیون U-PGx وجود دارند به دلیل پیچیدگی‌هایشان نبود. همه فنوتیپ‌ها بر اساس حضورشان در دستورالعمل‌ها و تعداد داروها با توصیه‌های عملی ارزیابی می‌شوند، جایی که عمل‌پذیر به‌عنوان «تعامل ژن-دارویی که نیاز به تغییر دارو، تنظیم دوز یا نظارت شدید دارد» تعریف می‌شود. برای تمام فارماکوژن‌های ذکر شده در دستورالعمل‌های CPIC و DPWG، تعداد تداخلات ژن-دارو محاسبه می‌شود.

### منبع

<https://www.nature.com/articles/s41397-021-00259-z>

PGx و ویژگی‌های بین ژنی. جایی که یک ویژگی به عنوان یک منطقه ژنومی مشروح شده مانند ژن‌های کدکننده پروتئین، مناطق تکراری قطعه‌ای، شبه‌ژن‌ها، و غیره تعریف می‌شود. تفسیر مرجع ژن کد برای ویژگی‌های ژنتیکی در ژنوم انسان (نسخه ۲۸) برای بررسی تفکیک‌پذیری هاپلوبلاک لوکوس‌های مهمی مانند ژن‌های کدکننده پروتئین پروژه Gencode با هدف طبقه‌بندی و شناسایی تمام ویژگی‌های ژن در ژنوم انسان از جمله تمام تفسیرها انجام می‌شود. برای هر ژنکد اتوزومال و ویژگی PGx، درصد ویژگی پوشش داده شده در هاپلوبلاک محاسبه می‌شود (تعداد جفت‌های پایه در هاپلو بلوک/طول ویژگی کل). مناطق با پوشش هاپلوبلاک ۹۰٪ به عنوان کاملاً فازی طبقه بندی می‌شوند، در حالی که مناطق بدون هاپلوبلاک همپوشانی بدون فاز باقی می‌مانند. همه مناطق دیگر به عنوان نیمه فازی مشخص شده‌اند.

تکرارهای قطعه‌ای (SD) و ترک‌های تکراری از مرورگر ژنوم UCSC (دانشگاه کالیفرنیا سانتا کروز) به دست می‌آیند. Bedtools برای شناسایی مناطق همپوشانی بین تمام ترک‌ها و فایل‌های تفسیر مورد بحث استفاده شد. برای هر لوکوس، درصد بخش‌هایی که با SD یا تکرار