

Od 2008 roku mo liwy jest odczyt sekwencji DNA (genomu) oraz RNA (transkryptomu) dzi ki technologii sekwencjonowania w ogólnie ci nazywanej Sekwencjonowaniem Nowej Generacji (NGS – Next Generation Sequencing). Wcze niej odbywało si to na mniejsz skal , metod Sangera u ywan od 1977 roku. Metoda ta była u ywana do rozszyfrowania kompletnego ludzkiego genomu (Human Genome Project) w roku 2000. W przeciwie stwie do metody Sangera, rozwi zania NGS generuj miliony krótkich odczytów (obecnie około 150 nukleotydowych par zasad) sekwencji zwanych odczytami (ang. reads). Nast pnie krótkie odczyty poddawane zostaj analizie, do której u ywa si specyficznych metod bioinformatycznych. Przykładem jest składanie w wi ksz sekwencje lub dopasowanie do istniejącego genomu.

Obecnie, pojedyncza maszyna dost pna na rynku, najbardziej zaawansowana technologicznie (Illumina 4000) mo e wygenerowa około 0.5T nukleotydowych par zasad w ci gu jednego dnia. Co powoduje, e odczyt ok. 12 całych ludzkich genomów mo e by uzyskany w 3.5 dnia pracy tej maszyny [1].

W analizie sekwencji istnieje wiele, ró nych typów/podej analizy statystycznej, które przetwarzaj informacje o sekwencjach DNA i RNA do postaci u ytecznej wiedzy dla nauk biologicznych b d medycyny. NGS we współpracy z technikami analizy, włączaj c przetwarzanie danych typu Big Data staje si podstaw dla nowego obszaru genomowej medycyny personalizowanej. Obszar ten został nazwany również „medycyną precyzyjną” w przemówieniu prezydenta Stanów Zjednoczonych Baracka Obamy z 30 stycznia 2015 roku, w którym ogłosił Inicjatywę Medycyny Precyzyjnej (ang. Precision Medicine Initiative). „Inicjatywa Medycyny Precyzyjnej ma wykorzystywa post py w dziedzinie genomiki, pojawiających si metod zarz dzania i analizy dużych zbiorów danych przy jednoczesnej ochronie prywatności, wykorzystaniu technologii informacji na temat zdrowia w celu przyspieszenia odkryć biomedycznych” [2].

Znaczenie tej inicjatywy oraz opartych na genomie rozwi zaniach medycznych w ogólnie ci przedstawiono w pozostałej części o wiadczenia Białego Domu:

“Medycyna Precyzyjna jest innowacyjnym podej ciem do zapobiegania chorobom oraz ich leczenia, która uwzgl dnia indywidualne ró nice w ludzkich genach, rodowisku oraz stylu ycia. Precyzyjna medycyna daje lekarzom klinicznym lepiej zrozumie zło one mechanizmy le ce u podstaw zdrowia pacjenta, choroby lub jej etapu. Pozwoli lepiej przewidzie , które zabiegi b d najbardziej skuteczne. Post py w medycynie precyzyjnej doprowadziły do nowych odkryć oraz wprowadzenia nowych metod leczenia, które s dostosowane do szczególnych cech fizycznych, takich jak profil genetyczny danej osoby czy indywidualne cechy nowotworu. To prowadzi do zupełnej zmiany sposobu leczenia chorób, takich jak rak. Przykładowo, pacjenci z dowolnym typem nowotworu b d poddawani rutynowym testom molekularnym w ramach opieki nad pacjentem. Testy te umo liwi lekarzom wybra odpowiednie leki, które zwi ksz szans na prze ycie i zmniejsz ryzyko nara enia na niekorzystne skutki. ”

Obecnie cena sekwencjonowania spada: 1000 dolarów za genom od ko ca 2014 roku, podczas gdy jako ro nie ze wzgl du na bardziej stabilne i sprawdzone protokoły laboratoryjne. W skim gardłem NGS jest obecnie proces uzyskania u ytecznej wiedzy biologicznej z dużych zbiorów danych sekwencjonowania. Odpowiedzi na potrzeby precyzyjnej analizy dużych danych genomicznych jest dalszy rozwój w dziedzinie okre lanej od niedawna jako “data science” (nauka o danych) i w dziedzinie oblicze chmurowych. W ci gu ostatnich lat, paradygmaty rozwoju oprogramowania i przetwarzania danych zwane map-reduce framework, umo liwiaj dystrybucj przetwarzania na wirtualnym klastrze maszyn i zapewni skalowalno wyników w stosunku do rozmiaru danych. W ostatnich 2 latach tak e rozwi zania maszynowego uczenia stały si skalowalne, co pozwala im działa na setki i tysie ce komputerów równoległych w razie potrzeby. To otwiera drog do rozwoju nowych, specjalistycznych metod uczenia maszynowego, które b d w stanie rozwi za nowe klasy problemów analizy danych NGS typu Big Data.

Dotychczas stosowane techniki analizy danych NGS w wi kszoci pochodz z obszaru statystyki. W moim projekcie gEMiLia, chciałabym poć czy moje do wiadczenie w analizie danych genomowych, które zdobyłam podczas przygotowywania pracy doktorskiej

z nowoczesnymi rozwi zaniami maszynowego uczenia, które s jednym z nurtów bada Zakładu Systemów Informatycznych, w którym pracuj od 01.01.2004.

NGS oraz obliczenia chmurowe znalazły si w 12 „przełomowych technologii” w raporcie McKinsey’a, opublikowanym w 2013 roku [3]. Wskazuje on na 12 technologii, które mog wpłyn na powa ne zmiany gospodarcze w najbliż szych latach. W moim projekcie spodziewam si poć czenia ich w cało oraz dodania dodatkowego wa nego czynnika, którym jest u ycie metod z obszaru maszynowego uczenia. Obszar ten jest powszechnie znanym, ale ci gle ywym i u ytecznym. Takie poć czenie jest nowatorskie i dotychczas niewiele na ten temat zrobiono w wiatowej społeczno ci naukowej. Poć czenie genomowej wiedzy bioinformatycznej, oblicze chmurowych oraz maszynowego uczenia jest rzadko ci . Wierz , e jest wiele do osi gni cia poprzez uzyskanie efektu synergii z nich, z udziałem mojej codziennej pracy badawczej.

Mój plan skupia si na trzech cie kach analizy danych: jako ci danych NGS, analizy wariantów DNA oraz analizy RNA, które technologicznie maj wiele wspólnego. Wyniki b d miały głównie form nowych rozwi za algorytmicznych, które mo na uruchomi równolegle na wielu komputerach w rodowisku chmurowym. To w konsekwencji pozwala rozwi za nowe klasy projektów analizy danych w dziedzinie biologii molekularnej i medycyny personalizowanej.

1. Illumina 4000 specifications. <http://www.illumina.com/systems/hiseq-3000-4000/specifications.html>
2. FACT SHEET: President Obama’s Precision Medicine Initiative. Illumina 4000 specifications. <http://www.illumina.com/systems/hiseq-3000-4000/specifications.html>
3. Disruptive technologies: Advances that will transform life, business, and the global economy. http://www.mckinsey.com/insights/business_technology/disruptive_technologies