

NIH/dbGAP の動向

国立研究開発法人科学技術振興機構  
 バイオサイエンスデータベースセンター

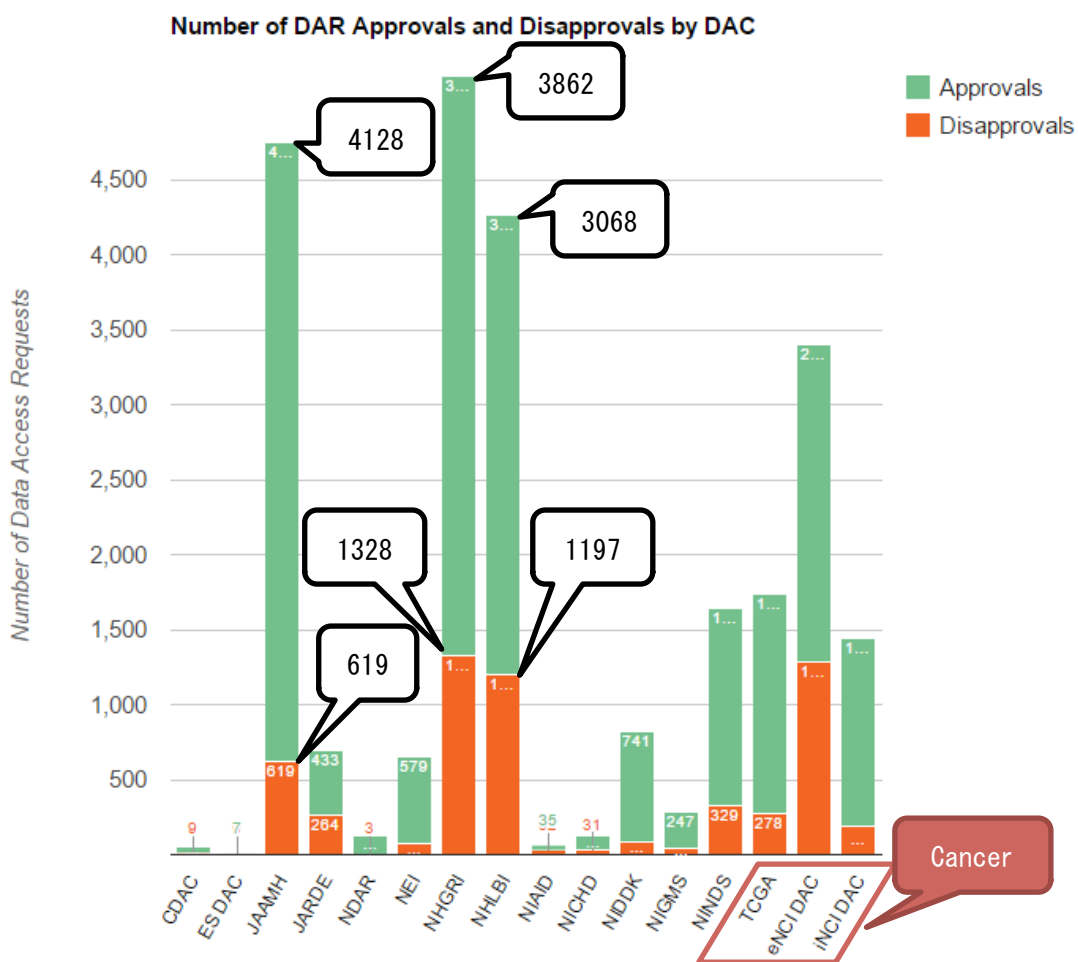
2015年5月19-21日にNCBIにて開催されたInternational Nucleotide Sequence Database Collaboration (INSDC)三極会議にDDBJの真島さん、児玉さん等が参加し、5月22日にdbGAPを訪問し、Dr. Michael Feolo、Dr. Lon Phan等と打合せをした。その報告を受け、特にData Access Committeeに関与する部分について報告する。

1：データ利用申請数および承認数（2007年以降）

利用申請数：30374、承認数：20272

2/3程度が承認されている。データの『利用制限（制限事項）』に抵触することから否認されている。

([http://gds.nih.gov/20ComplianceStatistics\\_dbGap.html#FN\\_violation\\_01](http://gds.nih.gov/20ComplianceStatistics_dbGap.html#FN_violation_01) のリスト参照)



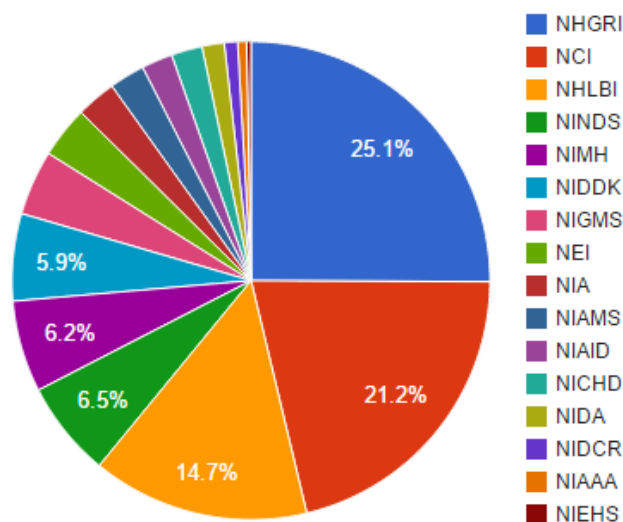
[http://gds.nih.gov/19dataaccesscommitteereview\\_dbGaP.html](http://gds.nih.gov/19dataaccesscommitteereview_dbGaP.html)

NIHには16のDACが存在する。それぞれのInstituteやCenterにDACが存在することが多い。JAAMHの様に薬物中毒・アルコール中毒・エイジング・メンタルヘルス4機関で1つのDACを所有していたり、The Cancer Genome Atlas (TCGA)の様にプロジェクトでDACを所有することもある。

2 : dbGAP に登録されているデータは、(1) National Human Genome Research Institute (NHGRI)、(2) National Cancer Institute (NCI)、(3) National Heart, Lung, and Blood Institute (NHLBI)の順で多い。341/350 研究 (97.4%) は米国からの登録 (日本からは1 研究)。

### NIH Institutes and Centers Sponsoring<sup>2</sup> dbGap Studies

Percentage of Sponsorship



[http://gds.nih.gov/17summary\\_dbGaP\\_statistics.html](http://gds.nih.gov/17summary_dbGaP_statistics.html)

3 :すでにデータの2次利用で950 報以上の論文が Publish されている(Data Use Report の記載を参照)。  
 >日本でも2次利用の論文を成果リストとしてNBDC ヒトDB から公開していきたい。

4 :100 万人以上の被験者のデータが登録されている。データキュレーションを実施し、同一検体からのデータは同じ内部ID が付けられている。新たに検出された変異はdbSNP やdbVar に登録している。  
 >NBDC ヒトDB でも、同一検体からのデータかどうか等のデータキュレーションを実施していきたい。

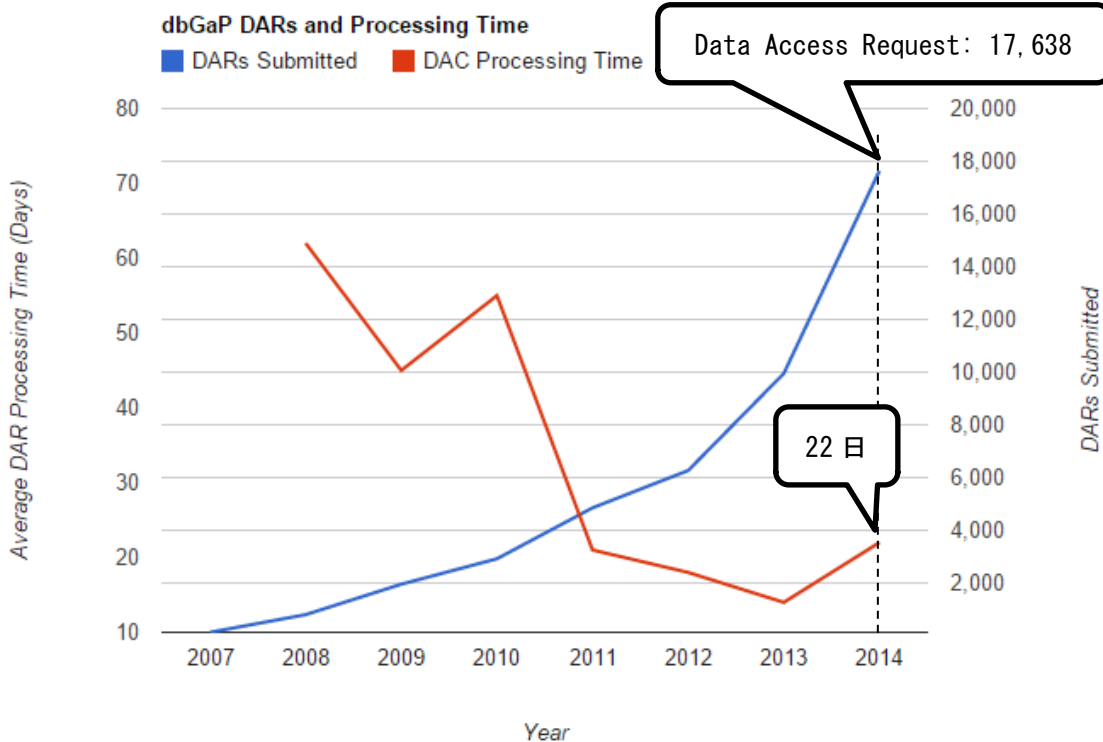
5 :GDS policy では、(1) NIH がFunding した全研究データをタイムリーに公開する(2) ヒトデータについては、Quality Control が終了次第、データを Submission し、6 か月以内にデータをリリースする(3) Publication Embargo 期間については、NIH による監視が難しく違反が起こりえるので、初めから設けないこととする、という運営を目指す。

6 :NBDC ヒトDB ではアクセス制限の表記をNIH のGenome-wide Association Policy に合わせて、『Open Data』と『Controlled Access Data』と名付けていたが、NIH のGDS に従い、NBDC でもOpen data をUnrestricted Access Data と名称変更することとする。

7 :Controlled access data を使用する際は、Genomic Data User Code of Conduct を遵守する必要が

ある。Genomic Data User Code of Conduct は『再配布の禁止』といった基本的項目が記載されていて（7項目）、NBDC ヒトデータ共有ガイドラインの中に記載済みの内容であった。

8 : データ利用申請に係る審査期間は平均 22 日



[http://gds.nih.gov/19timeprocessing\\_dbGaP.html](http://gds.nih.gov/19timeprocessing_dbGaP.html)

9 : データの質やデータ共有方法について NIH 基準を満たす組織は、Trusted Partner になることができる。NIH の Institute や Center と契約をすることで Trusted Partner になれる。ただ Data Access に関する審査は NIH の Institute や Center の DAC が実施する。

既に Trusted Partner となった組織 : Cancer Genomic Hub, Bionimbus, Cancer Genomics Data Commons, Cancer Genomics Cloud Pilot Project の 4 つ。

10 : 2015 年 3 月 23 日、クラウドサービスを利用した制限公開データの利用（解析や保管）に関する GDS policy への記載が公表された。クラウドサービス提供者ではなく、Controlled access data をクラウド環境で利用する者が所属する組織が責任を持つ、とした。データ利用時には、PI、組織の長、組織の IT 管理者がサインし、セキュアなクラウドを利用しているかどうかの責任は IT 管理者にある。クラウド環境によるデータ利用を希望する際は、利用申請時にクラウドを使用する旨を表明し、研究過程でどのような使用を予定しているか、また、どのプロバイダーのサービスを利用するのかについても報告しなければならない。現在、29 件利用している（大学・研究機関：17、民間：9、非営利：1、NIH：2）。プロバイダーは以下の通り。Amazon：11、Google：5、Private：5、DNAnexus：4、Microsoft：2、Seven Bridges：2

11 : データ利用申請時に DAC が確認していることは『どんな研究に使うか』と『Data Use Restriction』

が一致しているかどうか。Data Use Restrictionについては参考資料(Standard Data Use Limitations. pdf)参照

1 2 : データ利用者数 (PI 数) 3597 PIs (41 countries)

