# Cancer Data Sharing
# Genomic Data Commons

*Warren Kibbe, PhD*

*NCI Center for Biomedical Informatics*

@wakibbe

**NIH》 NATIONAL CANCER INSTITUTE**

The views expressed are my views and not those of NCI

June 21st, 2016

The Genomic Data Commons (GDC) is open!

Vice President Biden visited the GDC and spoke at ASCO on the 6[th] of June

**The New York Times**

Biden Unveiling Public Database for Clinical Data on Cancer

**FOX NEWS Health**

Biden unveiling public database for clinical data on cancer

**HealthData Management**

NCI launches open access resource to spur cancer research

**HOUSTON CHRONICLE**

Biden unveils searchable government cancer database

**SCIENTIFIC AMERICAN** **REUTERS**
PUBLIC HEALTH

Biden Unveils Major Database to Advance Cancer Research

**THE HUFFINGTON POST**
INFORM · INSPIRE · ENTERTAIN · EMPOWER

Biden Announces Crucial Piece Of His Cancer Moonshot Initiative

**The Washington Post**

Biden unveils launch of major, open-access database to advance cancer research

**FORTUNE**

Joe Biden Just Announced a Huge New National Cancer Database

Biden announces U.S. project to promote cancer data sharing

**REUTERS**

**CHICAGO SUN·TIMES**

VP Joe Biden in Chicago to promote Moonshot Initiative vs. cancer

**THE CANCER LETTER**

Biden Designates NCI's Genomic Data Commons As Foundation of Cancer Moonshot

**Daily Mail.com**

New US data system to centralize cancer information

**genomeweb**

NCI Launches Genomic Data Commons for Cancer Data Sharing

**HealthITAnalytics**
News and Resources for Healthcare Analytics Pros

NIH Launches Genomic Data Commons Supporting Cancer Moonshot

**fedscoop**

Biden launches data portal to back Cancer Moonshot
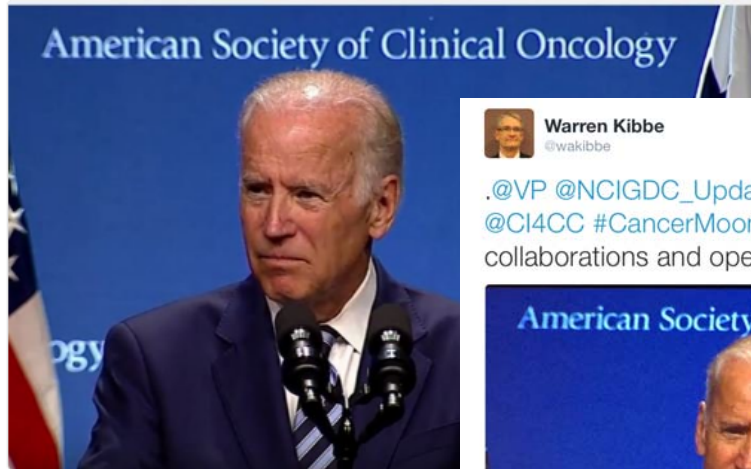
Shannon P. Hatch @sphatch · Jun 6
Colleagues from @UChicagoMed host @VP & @NCIDrDoug to launch #NCIGDC. #CancerMoonshot #BigData @theNCI
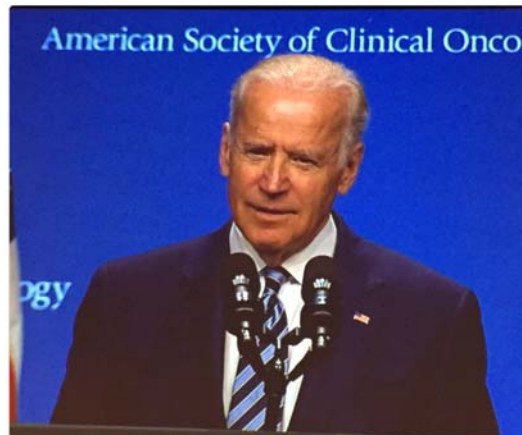
National Cancer Inst @theNCI · Jun 6
See @VP speak today at #ASCO2016 on the #CancerMoonshot
youtube.com/watch?v=s44YVo... #NCIGDC

American Society of Clinical Oncology

↩    ↻ 24    ♡ 22    •••

Warren Kibbe
@wakibbe

.@VP @NCIGDC_Updates #asco2016 @theNCI
@CI4CC #CancerMoonshot Need more
collaborations and open data sharing

American Society of Clinical Oncol

Warren Kibbe
@wakibbe

@VP @NCIGDC_Updates @theNCI @CI4CC
#CancerMoonshot #asco2016 genomic data
commons @NCIDrDoug

# Genomic Data Commons

The Cancer Genomic Data Commons (**GDC**) is an existing effort to standardize and simplify submission of genomic data to NCI and follow the principles of **FAIR** – Findable, Accessible, Interoperable, Reusable.

The GDC is part of the NIH Big Data to Knowledge (**BD2K**) initiative and an example of the **NIH Commons**

*Microattribution, nanopublications, tracking the use of data, annotation of data, use of algorithms, supports the data /software /metadata life cycle to provide credit and analyze impact of data, software, analytics, algorithm, curation and knowledge sharing*

# NCI Genomic Data Commons

- The GDC will go live with approximately 4.1 PB of data.

- This includes: 2.6 PB of legacy data;

- and 1.5 PB of "harmonized" data.

- 577,878 files about 14194 cases (patients), in 42 cancer types, across 29 primary sites.

- 10 major data types, ranging from Raw Sequencing Data, Raw Microarray Data, to Copy Number Variation, Simple Nucleotide Variation and Gene Expression.

- Data are derived from 17 different experimental strategies, with the major ones being RNA-Seq, WXS, WGS, miRNA-Seq, Genotyping Array and Expression Array.

# Genomic Data Commons Data Portal

# The NCI Genomic Data Commons User Interface
## Data Submission Dashboard

# Support the Precision Medicine Initiative
*The Genomic Data Commons and Cloud Pilots*

- Expand data model to include other data (e.g. imaging and proteomics)

- Allow easy publication of persistent links to data, annotations, algorithms, tools, workflows

- Measure usage and impact

- Change incentives for public contributions

# Cancer data ecosystem

| Discovery | Patient engaged Research | Surveillance Big Data Implementation research |
|---|---|---|
| Genomic research Clinical trials with genomic data | Clinical Research Observational studies | EHR, lab data, PROs, smart devices, decision support |
| Well characterized genomic data sets | Cancer cohorts | Patient data |



| Genomic information donor | Active research participation | SEER Learning from every cancer patient |
|---|---|---|