# The Genomic Data Analysis Network (GDAN)

**Lou Staudt**

**June 24, 2015**

# Computational Genomics –
# A Growing Necessity in Cancer Research

- TCGA production:
  - 33 tumor types and 11,500 cases
  - 2.5 petabytes (PB) of data
- Successful analysis and utilization of TCGA data required:
  - Experiments performed utilizing strict standardized protocols
  - Data in structured formats and available in public databases
  - Formation of Analysis Working Groups, with expertise in computational genomics, tumor biology and clinical oncology
- Genome Data Analysis Centers (GDACs) have been indispensible for progress in TCGA

# Genome Data Analysis Centers (GDACs)

## Generation of bioinformatics tools for the research community

# Firehose: An Automated Pipeline



http://gdac.broadinstitute.org

# cBio Cancer Genomics Portal

http://www.cbioportal.org/public-portal/



I'm interested in "X" gene pathway in colorectal cancer…



**Are there survival differences?**

**Aberrations in a specific genomic region?**

**I'm interested in patient TCGA-XXX**

# Genome Data Analysis Centers (GDACs)

## Data analysis for Analysis Working Groups

## Generation of clinically meaningful molecular subgroups of cancer

# Four Molecular Subgroups of Endometrial Cancer Defined by Integrative Analysis

# Molecular Subgroups Refine Histological Diagnosis Of Endometrial Carcinoma



POLE (ultra-mutated)

MSI (hypermutated)

Copy-number low (endometriod)

Copy-number high (serous-like)

Mutations Per Mb

PolE
MSI / MSH2
Copy #
PTEN
p53

Histology

n = 17    n = 65    n = 90    n = 60

Spectrum* (232)

POLE (ultramutated) (17)

MSI high (215)

MSI (hypermutated) (65)

(150)   CN cluster 4

Copy-number low (endometrioid) (90)    Copy-number high (serous-like) (60)

*(%[CA] > 0.2) AND (%[CG] < 0.03) AND (SNV count > 500)

Histology

Endometrioid

Serous

Serous misdiagnosed as endometrioid?

TCGA Nature 497:67 (2013)

# Molecular Diagnosis of Endometrial Cancer May Influence Choice of Therapy

# Integration Matters



12 tumor types

Leukemia (LAML)
Lung adenocarcinoma (LUAD)
Lung squamous (LUSC)
Kidney (KIRC)
Bladder (BLCA)
Endometrial (UCEC)

Glioblastoma (GBM)
Head and neck (HNSC)
Breast (BRCA)
Ovarian (OV)
Colon (COAD)
Rectum (READ)

Thematic pathways

Omics characterizations

Platforms

Mutation
Copy number
Gene expression
DNA methylation
MicroRNA
RPPA
Clinical data

BRCA
BLCA
COAD
GBM
HNSC
KIRC
LAML
LUAD
LUSC
OV
READ
UCEC

Samples

Genes/loci

# PanCan Analysis Reveals Clinically Distinct Bladder Cancer Subtypes



Hoadley et al.  Cell 2014 158:929-44

National Cancer Institute

# PanCan Analysis Reveals Clinically Distinct Bladder Cancer Subtypes

Hoadley et al.  Cell 2014 158:929-44

# PanCan Analysis Reveals Clinically Distinct Bladder Cancer Subtypes



Hoadley et al.  Cell 2014 158:929-44

National Cancer Institute

# Computational Genomics for Center for Cancer Genomics Initiatives

- CCG initiatives will:
  - Conduct comprehensive genome-wide analyses of molecular alterations in cancers
  - Utilize multiple platforms to profile the genome, transcriptome and epigenome of cancer
- CCG goals include:
  - Identify genomic alterations that influence the development of cancer and the response to treatment
  - Collaborate with other NCI Divisions and Centers to conduct the most meaningful genomic studies
  - Support the Precision Medicine Initiative

# The CCG Genomics Pipeline

## Cancer Biopsies

## Biospecimen Core Repository (BCR)

### Tumor Pathology QC
- % Tumor Nuclei
- % Necrosis
- Dx Confirmation via histology and pathology report

### Molecular Analyte QC
- Spectrophotometry
- RNA Bioanalyzer
- Electrophoresis
- Genotyping

## Genome Characterization Centers

Exome seq                RNA-seq
Whole genome seq         DNA Methylation

## Genome Data Analysis Network (GDAN)

| Genetic aberrations | Data analysis: | Data integration: |
|---|---|---|
| Mutations | Molecular subgroups | Functional vs. structural |
| Copy number | Co-occurrence / exclusion | Master regulator analysis |
| Translocations | Comparison to TCGA | Pathway analysis |

itute

# Projects Involving the GDAN

- CCG initiatives (some with other NCI Divisions):
  - Cancer Driver Discovery Program (CDDP)
  - The Adjuvant Lung Cancer Enrichment Marker Identification and Sequencing Trials (ALCHEMIST)
  - Exceptional Responders (in collaboration with DCTD)
  - Clinical Trials Sequencing Program (in collaboration with DCTD)
  - Environment and Genetics in Lung Cancer Etiology (EAGLE, in collaboration with DCEG)
- The GDAN can be used to support any NCI project that utilizes the CCG genomics pipeline

National Cancer Institute

# Composition of the GDAN

- **Processing GDAC**
  - Develops and implements appropriate bioinformatic systems for rapid high-throughput processing
  - Operates closely with the NCI Genomic Data Commons (GDC) to generate primary genomic results
  - One center will be awarded

- **Visualization GDACs**
  - Provides user-friendly bioinformatics tools and data portals for the exploration of results
  - Explores new methods to integrate data
  - Two centers will be awarded

- **Specialized GDACs**
  - Provides in-depth expertise on individual platforms
  - Provides analytical support to Analysis Working groups
  - Eleven centers will be awarded

National Cancer Institute

# Mechanisms of Award & Budget

- All awards will be U24 Cooperative Agreements
- Budget is as follows (in thousand dollars):

| GDAC Type | Award Number | Amount /Year | FY2016 | FY2017 | FY2018 | FY2019 | FY2020 |
|---|---|---|---|---|---|---|---|
| Process | 1 | 1,000 | 1,000 | 1,000 | 1,000 | 1,000 | 1,000 |
| Visual | 2 | 1,000 | 2,000 | 2,000 | 2,000 | 2,000 | 2,000 |
| Special | 11 | 500 | 5,500 | 5,500 | 5,500 | 5,500 | 5,500 |
|  |  | Total | 8,500 | 8,500 | 8,500 | 8,500 | 8,500 |
| Grand Total |  | 42,500 |  |  |  |  |  |

National Cancer Institute

# Justification for the GDAN RFA

- TCGA experience suggests that data analysis in large-scale genomic characterization programs requires a coordinated group of experts in computational genomics

- This coordinated network requires a detailed statement of needs, including time lines and deliverables

- It is unlikely that such a network would evolve from a disparate collection of investigator-initiated grants

- The GDAN will support and stimulate the development of computational genomics tools and methodologies for the research community

National Cancer Institute

# Justification for Cooperative Agreement

- The CCG genomics pipeline requires coordination of:
  - Biospecimen processing
  - Genomic characterization of analytes
  - Analysis of the resulting data
- This coordination is maintained by the CCG Program staff working with the Analysis Working Groups.
- A cooperative agreement will allow CCG Program staff to deploy GDAN centers strategically to meet NCI needs
- A cooperative agreement will ensure that all results will be made publically available on a defined timeline
- The cooperative agreement will require that all bioinformatics tools be open-source and publically available

National Cancer Institute

# Evaluation Criteria

- The impact of the GDAN will be judged by:
  - Successful and timely support of the Analysis Working Groups (AWGs) for each CCG/NCI project
  - Cancer relevance of publications supported by the GDAN, as measured by citations and other metrics
  - Adoption of the bioinformatics tools generated by the GDAN for data processing and visualization
  - Training and support of trainees in computational

National Cancer Institute

# Questions?

National Cancer Institute