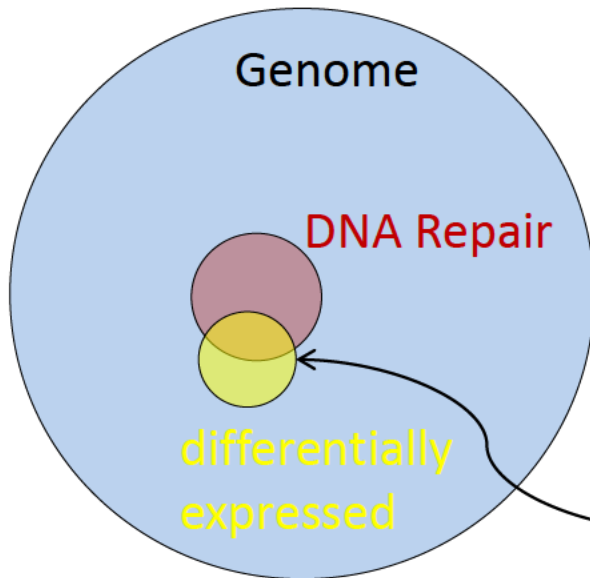# Outline

- Statistical significance for gene annotations
- Big data
  - L1000 transcriptional assay
  - Chemical sensitivity dataset
  - PubChem
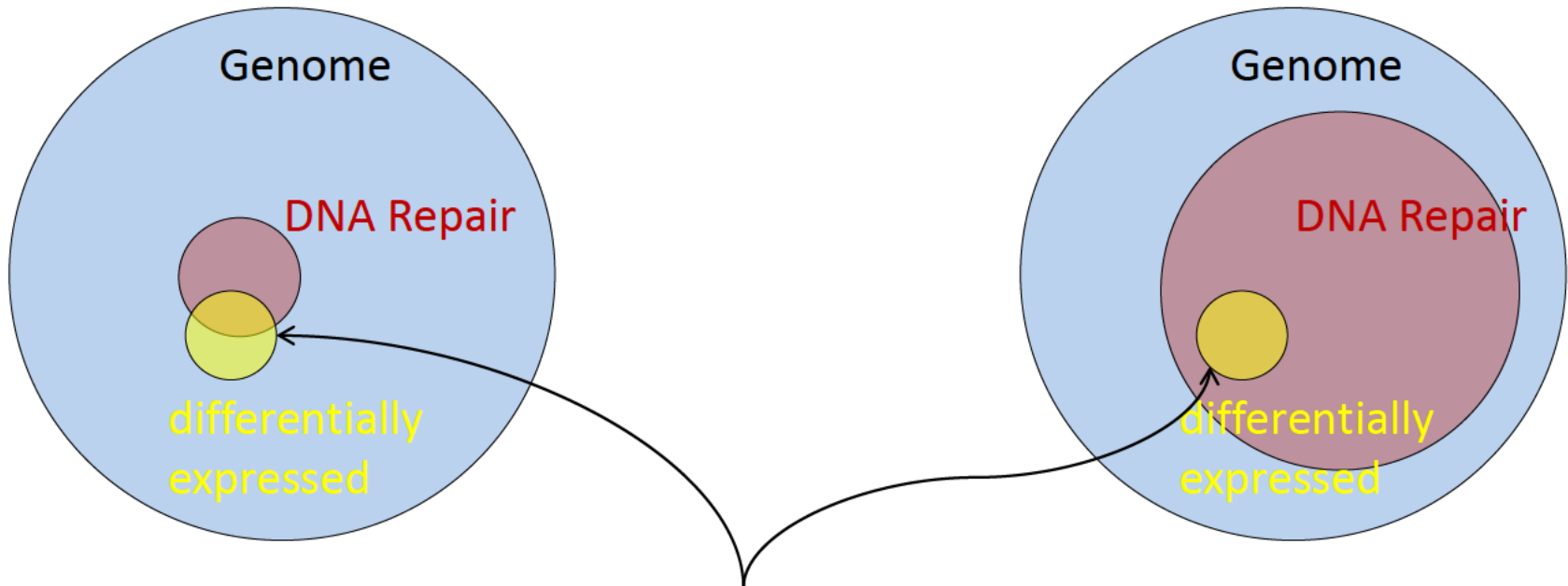  - TCGA
  - Drug Repositioning

# Statistical significance



I found that ten of the upregulated genes in my dataset are annotated as "DNA Repair" ...

Is this overlap significant?

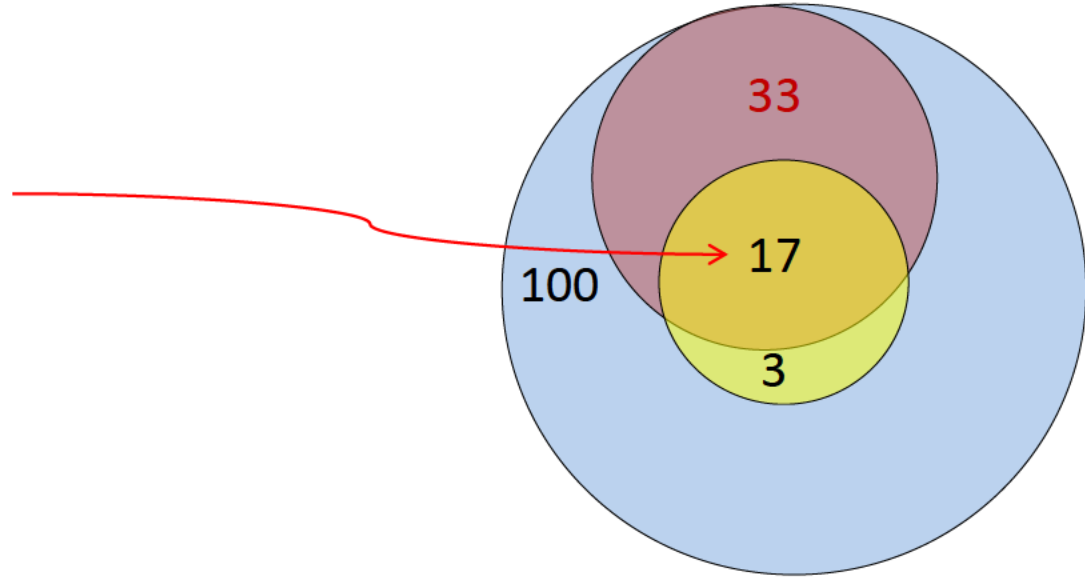To answer this question we need a null model.

# Statistical significance

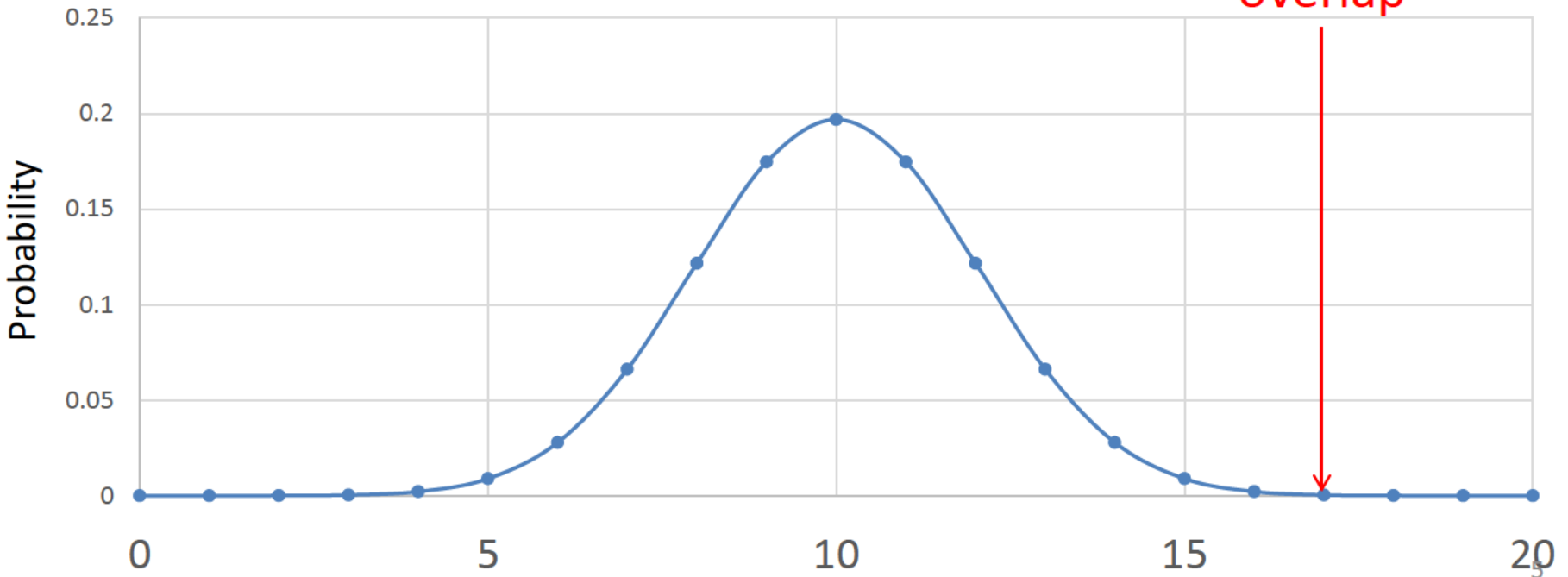The significance depends on the size of the lists.

Genome

DNA Repair

differentially expressed

Genome

DNA Repair

differentially expressed

If the two lists had nothing in common, could we still get this degree of overlap?
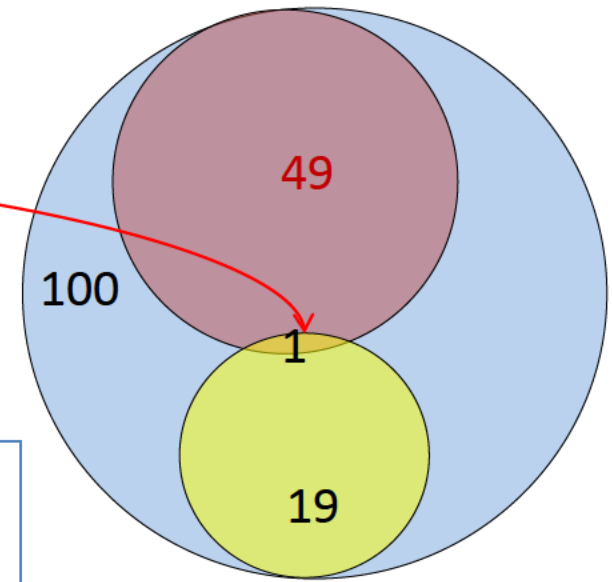
There are 17 overlapping genes. Is that surprising?

33

17

100

3

## Hypergeometric Distribution

Observed overlap

Probability

0    5    10    15    20

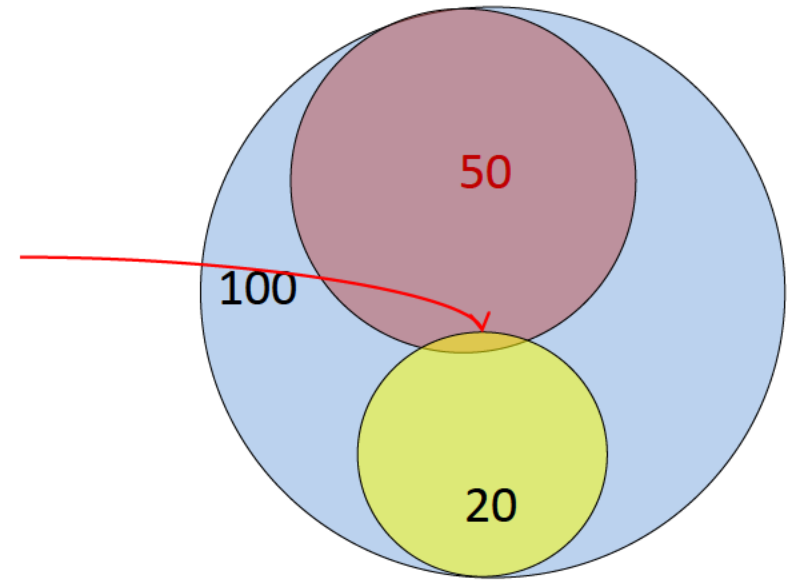There is only one overlapping gene.
Is that surprising?

49

100

1

19

Yes! You would expect to see a **_larger_** overlap under the null model.

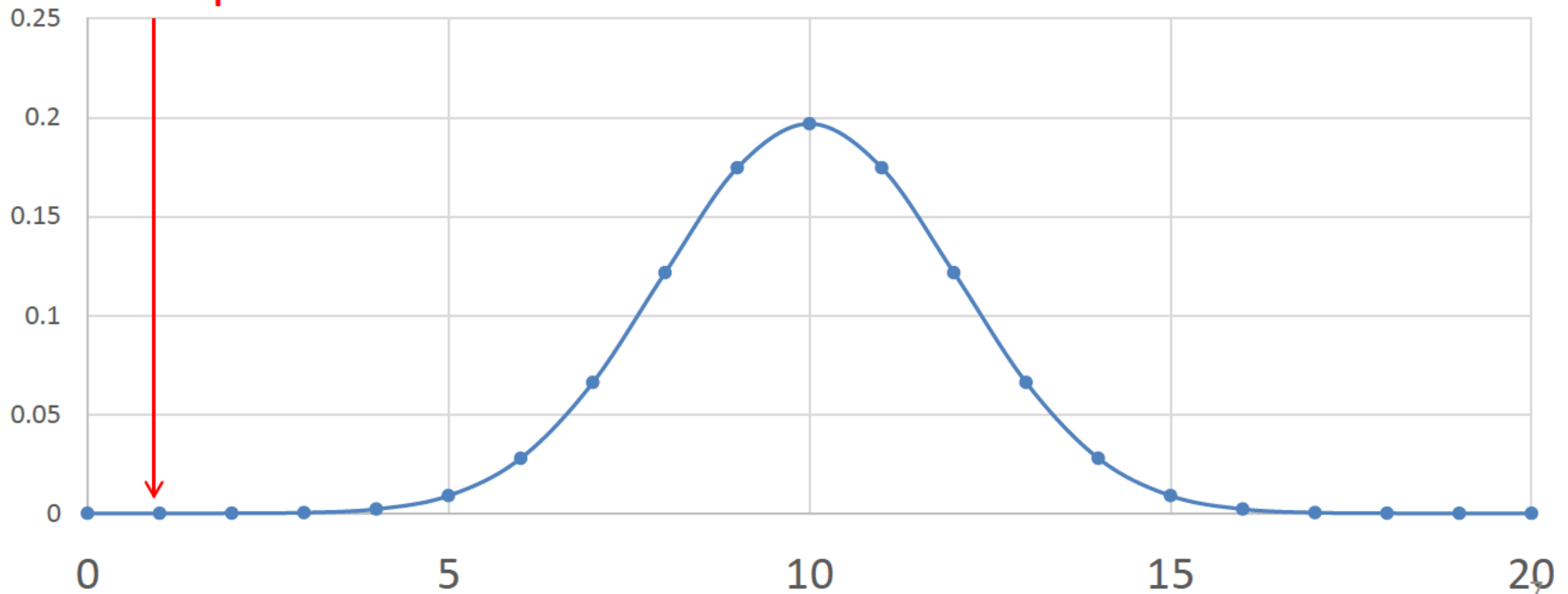Are the yellow genes **_enriched_** for the red function?

No! Quite the opposite!

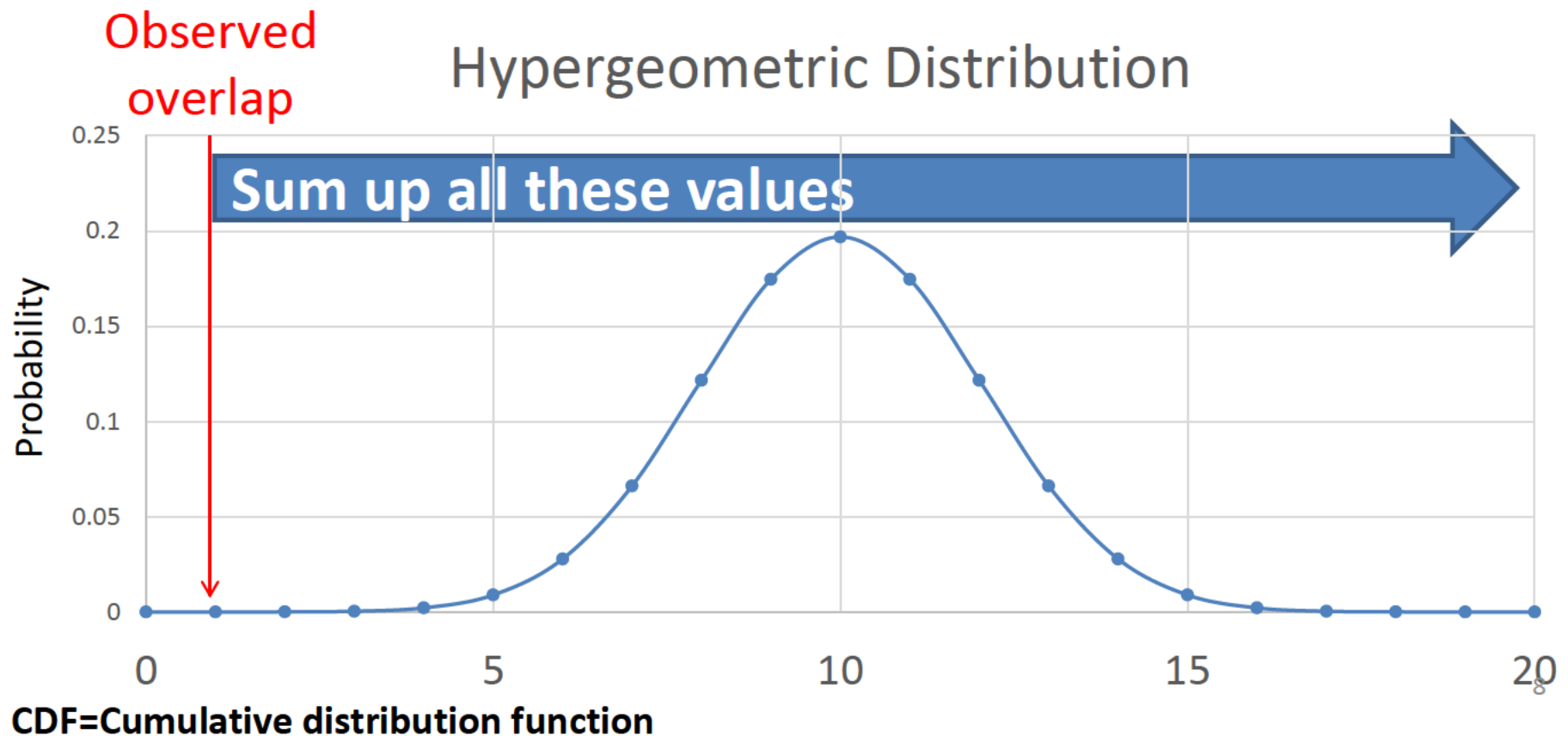The Hypergeometric p-value is the probability of observing an exact overlap

50

100

20

Observed overlap

Hypergeometric Distribution

Probability

# The CDF helps us find enriched terms

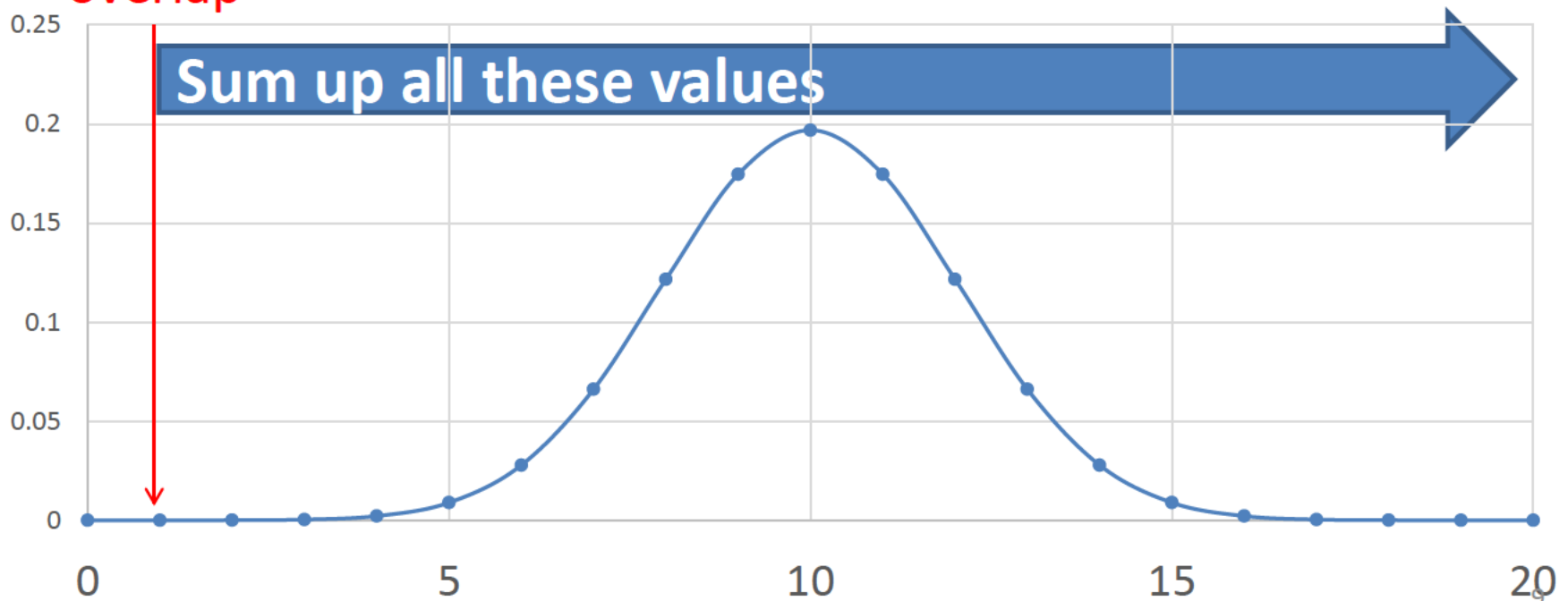We want to compute the probability of observing **at least** this overlap under our null model.

**CDF=Cumulative distribution function**

# The CDF helps us find enriched terms

$$CDF(Overlap) = \sum_{n=overlap}^{\substack{Number\ of \\ genes\ in\ DNA\ Repair}} \frac{\binom{DNA\ repair}{n}\binom{Genome - DNA\ repair}{DiffExp - n}}{\binom{Genome}{DiffExp}}$$
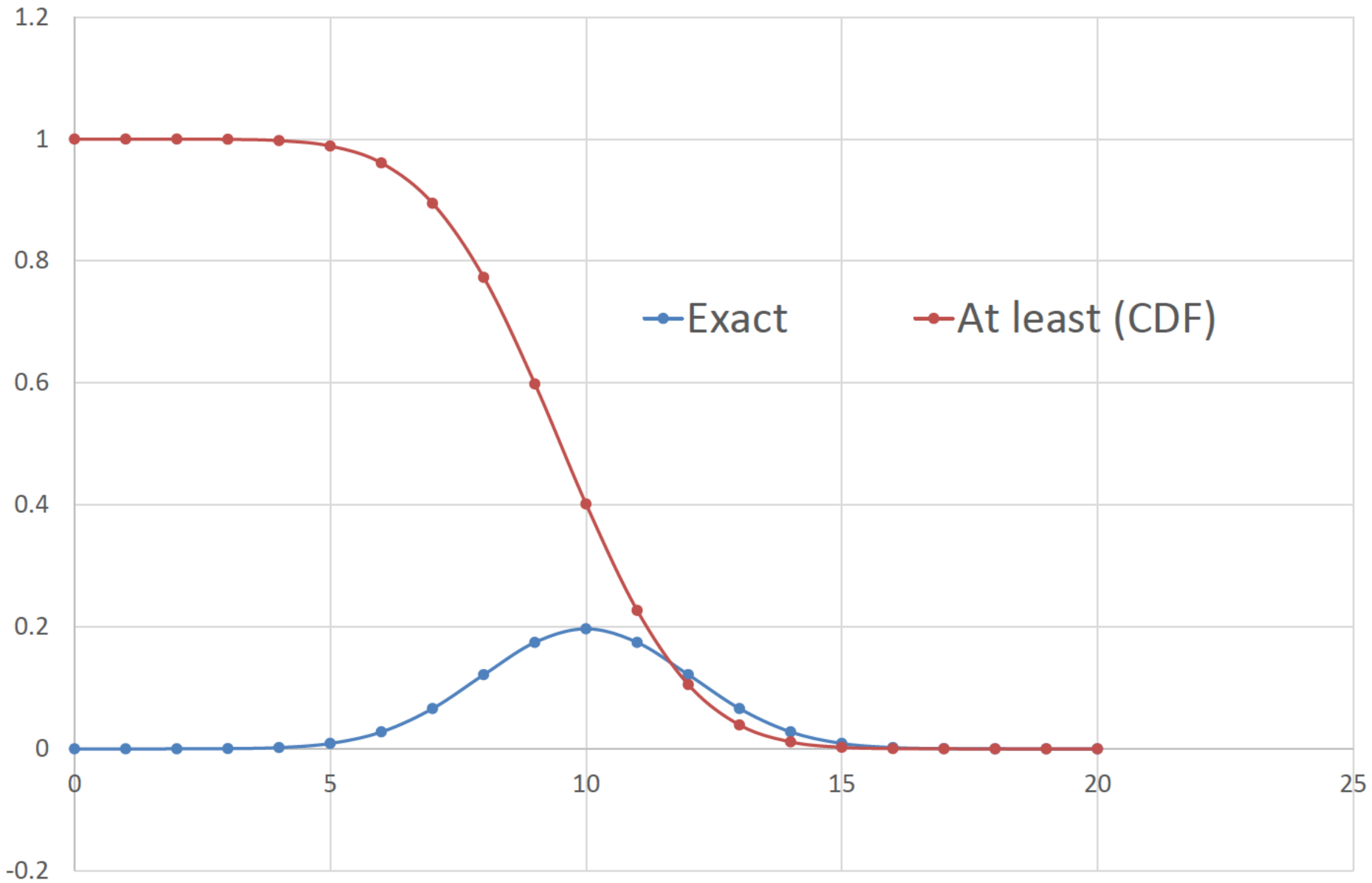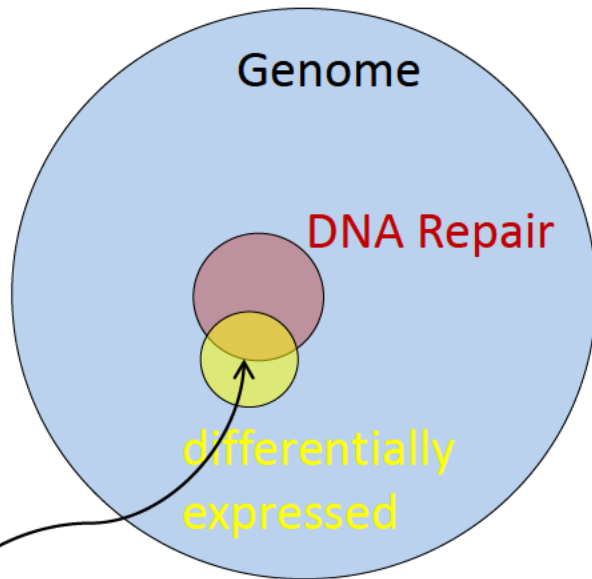
Hypergeometric Distribution

**Observed overlap**

**Sum up all these values**



CDF=Cumulative distribution function

# Hypergeometric

# Statistical significance



Genome

DNA Repair

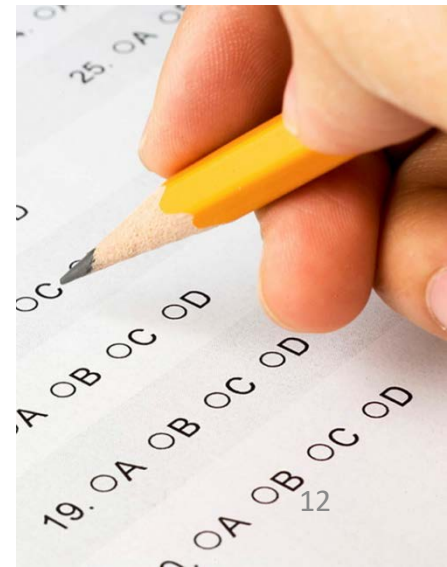differentially expressed

Is this overlap significant?

- We wish to test if a term is "enriched" in our data.
- But the hypergeometric gives the probability of getting **exactly** this amount of overlap for two randomly chosen sets of genes of the same size.
- Using the CDF, we can ask if we see **_more_** of a term than we would expect under the null model.

# Testing Multiple Hypotheses

- Example:  Filter GO terms using a p-value threshold of 0.01

- By definition, the null-hypothesis has a 1% probability of being correct ***for each test.***

- There are roughly 30,000 terms in GO.

- At this level, we expect roughly 300 false positives!
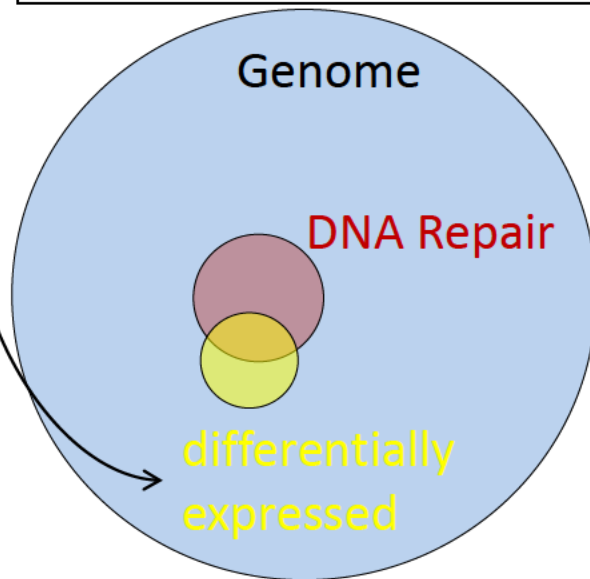
# Multiple Hypotheses

- A simple solution: require that the p-value be small enough to reduce the false positives to the desired level.

- This is called the Bonferroni correction.

- In our case, we would only accept terms with a

$$p \leq \frac{0.01}{30,000} = \frac{desired\ threshold}{number\ of\ tests}$$

- Since our tests are not all independent, this is very conservative, and will miss many true positives

- More sophisticated approaches exist, such as controlling the "false discovery rate".
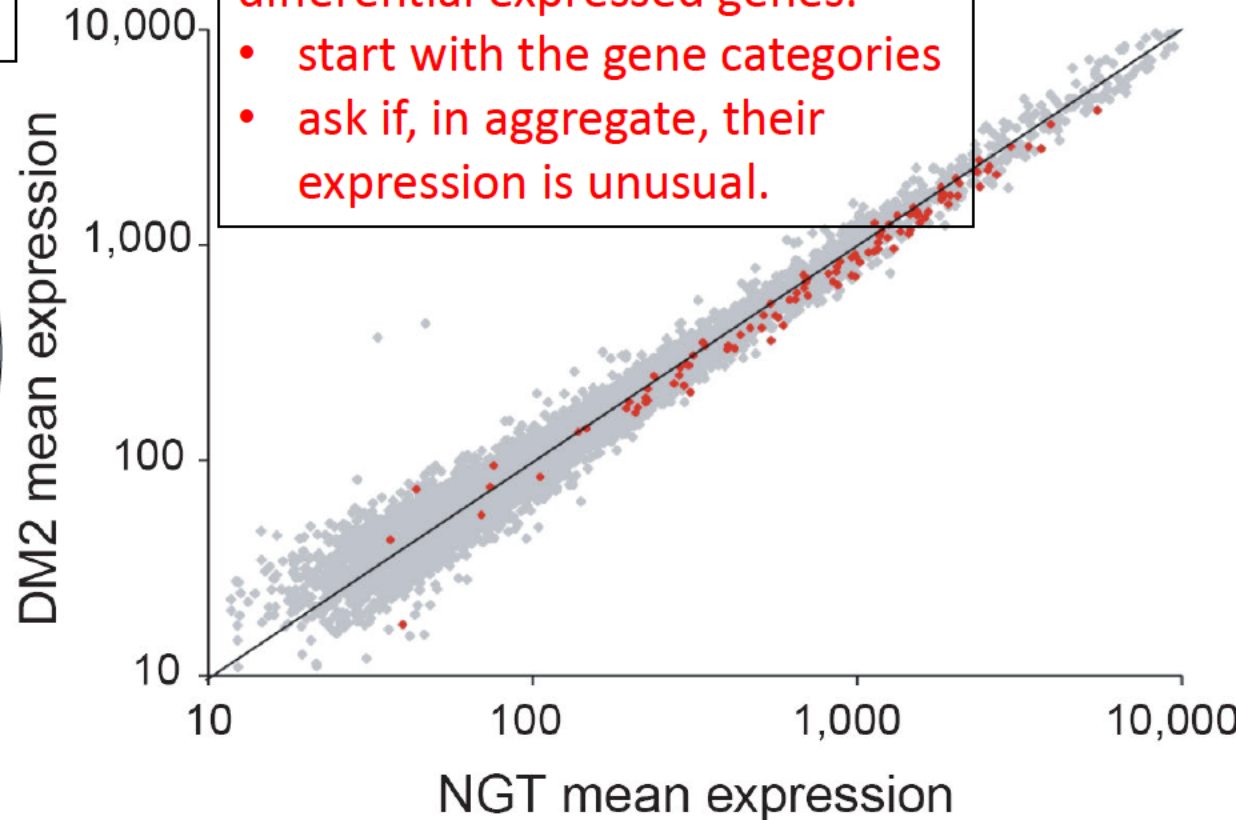
# Aggregate score statistics

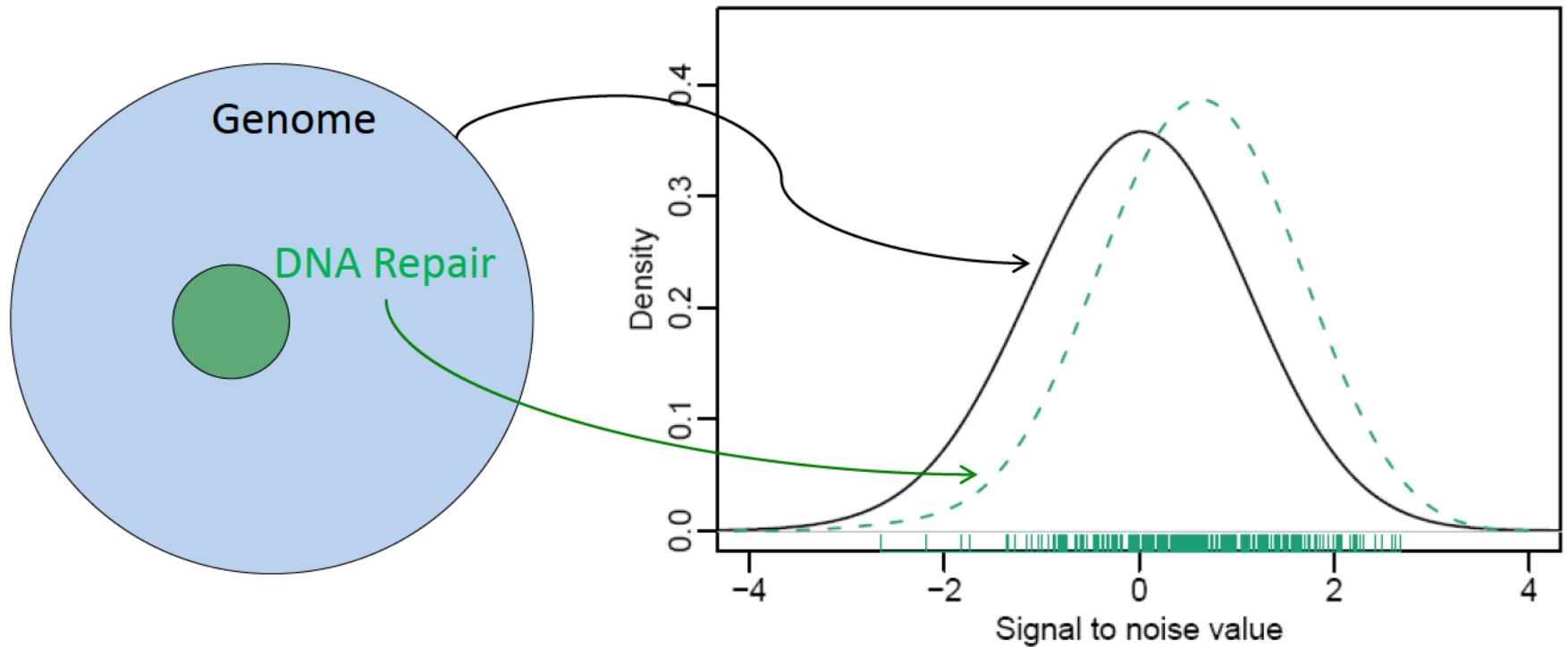My results depend on how I defined "differentially expressed"

Genome

DNA Repair

differentially expressed

Instead of starting with differential expressed genes:
- start with the gene categories
- ask if, in aggregate, their expression is unusual.



DM2 mean expression

NGT mean expression

Mootha *et al.* (2003). *Nature Genetics* **34**, 267 – 273. doi:10.1038/ng1180

# Aggregate score statistics

# Aggregate score statistics

http://www.broadinstitute.org/gsea/

# Learning Objectives

- To understand types and sources of biological "big data" and how they are used

# Big Data Creates an Opportunity

## Transcription

**Genomics**



**Analysis of protein-coding genetic variation in 60,706 humans**

**The Exome Aggregation Consortium**

**ExAC**

**>2.3 million samples so far**

**>94 million compounds**

**2.5 Petabytes**

**33 types of tumors**

**11,000 patients**

**7 data types**

18

Example 1

# L1000:  A VERY LARGE TRANSCRIPTIONAL DATASET

# Cell

# A Next Generation Connectivity Map: L1000 Platform and the First 1,000,000 Profiles

20

# Clustering Transcriptional Results

PCL=perturbagen class (targets identified)
Discovery = target unknown

C

Topoisomerase inhibitors

PIK3/MTOR inhibitors

HSP inhibitors

Protein synthesis inhibitors

ATPase inhibitors

HDAC inhibitors

NFkB/IKK inhibitors

Proteosome inhibitors

Aurora kinase inhibitors

HIF modulators

1 MEK inhibitors
2 Tubulin inhibitors
3 Glucocorticoid receptor agonists

Discover (2,418)
PCL member (625)

B

| Affinity to target | BRD-5657 kD (nM) | BRD-5161 kD (nM) | BRD-9186 kD (nM) |
|---|---|---|---|
| AKT1 | > 10,000 | > 10,000 | > 10,000 |
| MTOR | 87 | 1,900 | 2,600 |
| PIK3CA | 680 | 95 | 7,200 |
| PIK3CB | 1,000 | 1,200 | > 10,000 |
| PIK3CD | 1,200 | 480 | > 10,000 |
| PIK3CG | 330 | 46 | > 10,000 |

PIK3/MTOR inhibitors

Example 2

# A VERY LARGE SENSITIVITY ASSAY

# Harnessing Connectivity in a Large-Scale Small-Molecule Sensitivity Dataset.

481 compounds, including FDA-approved drugs, clinical candidates, and small-molecule probes



Seashore-Ludlow *et al.*
*Cancer Discovery (2015) 5(11); 1210–23*

25

The 481 compounds were tested at 16 concentrations in duplicate against 664 cancer cell lines.

# Cluster often represent common sensitivity to a mechanism

# Discovery of new way to target neuroblastoma?



IGF1R inhibitors

- NVP-TAE684
- BMS-536924
- BMS-754807
- NVP-ADW742
- Linsitinib

0.26

NVP-TAE684

| nmol/L | 0 | 0 | 10 | 30 | 90 | 270 | |
|--------|---|---|----|----|----|-----|---|
| IGF1 | – | + | + | + | + | + | |
| | | | | | | | ALK |
| | | | | | | | pALK |
| | | | | | | | IGF1R |
| | | | | | | | pIGF1R |
| | | | | | | | AKT |
| | | | | | | | pAKT |
| | | | | | | | GAPDH |

# Discovery of new way to target neuroblastoma?



Most cells show a similar response to NVP-TAE684 and BMS-754807, suggesting a similar mechanism

Cells with high levels of ALK (NB1) or specific ALK genomic changes are more sensitive to NVP-TAE684, which targets both IGF1R and ALK than they are to the IGF1R inhibitor BMS-754807

# NB1 responds to a combination of ALK and IGF1R inhibitors



Known ALK Inhibitor

Example 3

# PUBCHEM: A DATABASE OF CHEMICAL COMPOUNDS

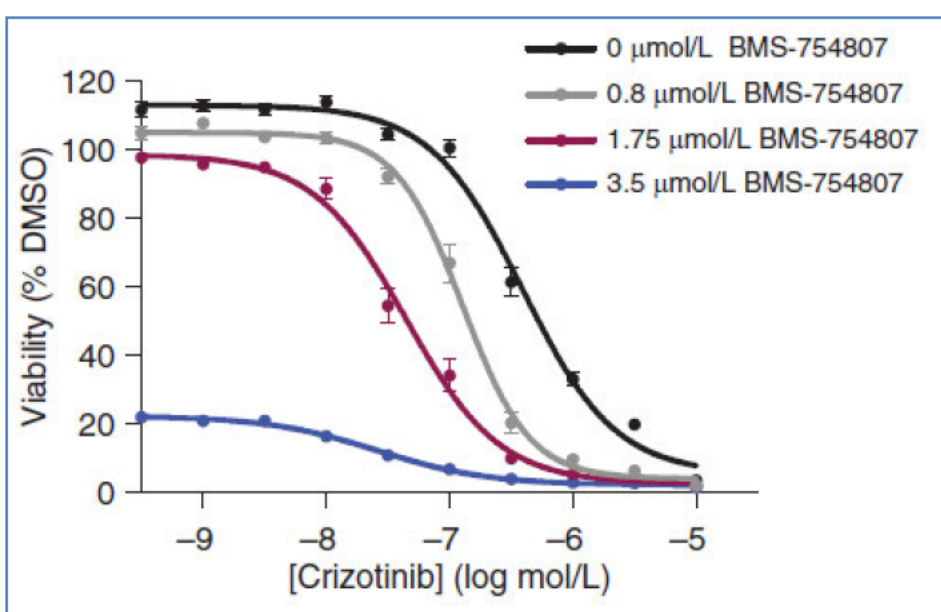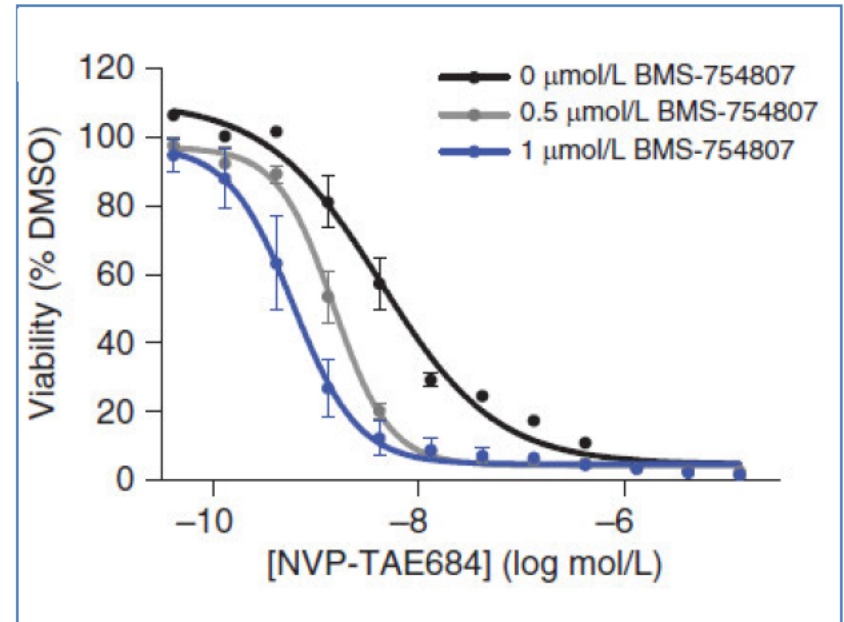Databases > | Upload | Services > | Help | more > | Today's Statistics >

| | |
|---|---|
| Compounds: | 94,703,715 |
| Substances: | 242,313,312 |
| BioAssays: | 1,252,878 |
| Tested Compounds: | 2,570,179 |
| Tested Substances: | 4,157,676 |
| RNAi BioAssays: | 170 |
| BioActivities: | 234,773,916 |
| Protein Targets: | 10,857 |
| Gene Targets: | 22,106 |

**Pub**

BioAssay ? | Com... | Substance

Go | Limits | Advanced

Try the PubChem Search Beta

New PubChem presents at the American Chemical Society National Meeting in New Orleans (March 18-22, 2018). Read more...

more ...

BioAssay Tools

Structure Search

3D Conformer Tools

Structure Clustering

Classification

Upload

Download

PubChem FTP

Write to Helpdesk | Disclaimer | Privacy Statement | Accessibility | Data Citation Guidelines
National Center for Biotechnology Information
NLM | NIH | HHS

32

# Tae-684

▶ Cite this Record

STRUCTURE    VENDORS    LITERATURE    PATENTS    BIOACTIVITIES

| | |
|---|---|
| **PubChem CID:** | 16038120 |
| **Chemical Names:** | NVP-TAE684; 761439-42-3; NVP-TAE 684; TAE684; TAE-684; 5-chloro-N4-(2-(isopropylsulfonyl)phenyl)-N2-(2-methoxy-4-(4-(4-methylpiperazin-1-yl)piperidin-1-yl)phenyl)pyrimidine-2,4-diamine   More... |
| **Molecular Formula:** | $C_{30}H_{40}ClN_7O_3S$ |
| **Molecular Weight:** | 614.206 g/mol |
| **InChI Key:** | QQWUGDVOUVUTOY-UHFFFAOYSA-N |
| **Substance Registry:** | FDA UNII |

PUBCHEM > COMPOUND > TAE-684     *Modify Date: 2018-03-17; Create Date: 2007-04-09*

## Contents «

## 1 2D Structure

Q Search    ⬇ Download    🖼 Get Image



Q Magnify

▶ *from PubChem*

33

Total Pages: 96    Display: [20 ▾] [Go To] Page [1]    [◄◄] [◄] [►] [►►]

Sort: ○ ▲ ◉ ▼ [Click the result table header to sort]
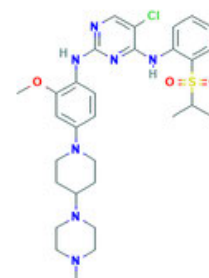
| # | ☐ | AID | Activity | AC≤1[µM] | AC≤1[nM]▾ | AC Range | BioAssay [Outcome Type] | Protein Target |
|---|---|---|---|---|---|---|---|---|
| 1 | ☐ | 624743 | ▬ | 1 | 1 | 0.00095 [µM] | Binding constant for LTK kinase domain [Confirmatory] | Leukocyte tyrosine kinase receptor[gi:143811416] |
| 2 | ☐ | 742112 | ▬ | 1 | 1 | 0.000927 [µM] | SANGER: Inhibition of human DEL cell growth in a cell viability assay. [Confirmatory] | |
| 3 | ☐ | 624825 | ▬ | 1 | 1 | 0.00085 [µM] | Binding constant for BMPR1B kinase domain [Confirmatory] | Bone morphogenetic protein receptor type-1B[gi:6226778] |
| 4 | ☐ | 624916 | ▬ | 1 | 1 | 0.00083 [µM] | Binding constant for ULK1 kinase domain [Confirmatory] | Serine/threonine-protein kinase ULK1[gi:317373288] |
| 5 | ☐ | 625076 | ▬ | 1 | 1 | 0.00093 [µM] | Binding constant for PLK4 kinase domain [Confirmatory] | Serine/threonine-protein kinase PLK4[gi:160113150] |
| 6 | ☐ | 741847 | ▬ | 1 | 1 | 6.03e-05 [µM] | SANGER: Inhibition of human SCC-3 cell growth in a cell viability assay. [Confirmatory] | |
| 7 | ☐ | 741855 | ▬ | 1 | 1 | 0.000564 [µM] | SANGER: Inhibition of human SF539 cell growth in a cell viability assay. [Confirmatory] | |
| 8 | ☐ | 624899 | ▬ | 1 | 1 | 0.00049 [µM] | Binding constant for ROS1 kinase domain [Confirmatory] | Proto-oncogene tyrosine-protein kinase ROS[gi:126302596] |
| 9 | ☐ | 624741 | ▬ | 1 | 1 | 0.00086 [µM] | Binding constant for LRRK2(G2019S) kinase domain [Confirmatory] | Leucine-rich repeat serine/threonine-protein kinase 2[gi:294862450] |
| 10 | ☐ | 624742 | 🗓 | | | 10 [µM] | Binding constant for NEK5 kinase domain [Confirmatory] | Serine/threonine-protein kinase Nek5[gi:74758252] |
| 11 | ☐ | 624900 | ▬ | 1 | | 0.16 [µM] | Binding constant for RSK1(Kin.Dom.1-N-terminal) kinase domain [Confirmatory] | Ribosomal protein S6 kinase alpha-1[gi:20178306] |
| 12 | ☐ | 624901 | ▬ | 1 | | 0.65 [µM] | Binding constant for RSK1(Kin.Dom.2-C-terminal) kinase domain [Confirmatory] | Ribosomal protein S6 kinase alpha-1[gi:20178306] |
| 13 | ☐ | 624902 | ▬ | 1 | | 0.15 [µM] | Binding constant for MEK4 kinase domain [Confirmatory] | Dual specificity mitogen-activated protein kinase kinase 4[gi:1170596] |
| 14 | ☐ | 624903 | ▬ | | | 4.8 [µM] | Binding constant for SRPK1 kinase domain [Confirmatory] | SRSF protein kinase 1[gi:209572680] |
| 15 | ☐ | 624904 | 🗓 | | | 10 [µM] | Binding constant for NEK4 kinase domain [Confirmatory] | Serine/threonine-protein kinase Nek4[gi:229462924] |

**Targets with Kd or IC50 below 10nM:**

ABL1 ALK BMPR1B DCLK1 EGFR FER FES FLT3 GAK IGF1R INSR INSRR LRRK2 LTK NUAK2 PLK4 PTK2 PTK2B ROS1 STK33 TNK1 TNK2 ULK1 ULK2 YES1

Example 4

# TCGA: MULTI-OMIC TUMOR DATA

**What questions might you ask using these sequencing data?**

# Top Mutated Cancer Genes

# ⊞ Summary

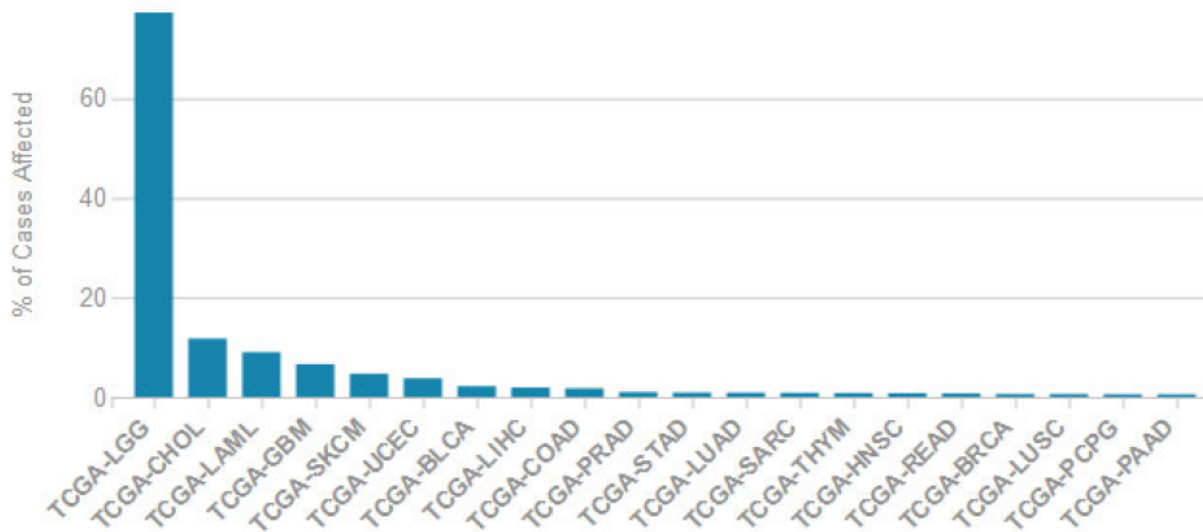| | |
|---|---|
| **Symbol** | IDH1 |
| **Name** | isocitrate dehydrogenase 1 (NADP+), soluble |
| **Synonyms** | -- |
| **Type** | protein_coding |
| **Location** | chr2:208236227-208266074 (GRCh38) |
| **Strand** | — |

**Description** Isocitrate dehydrogenases catalyze the oxidative decarboxylation of isocitrate to 2-oxoglutarate. These enzymes belong to two distinct subclasses, one of which utilizes NAD(+) as the electron acceptor and the other NADP(+). Five isocitrate dehydrogen...
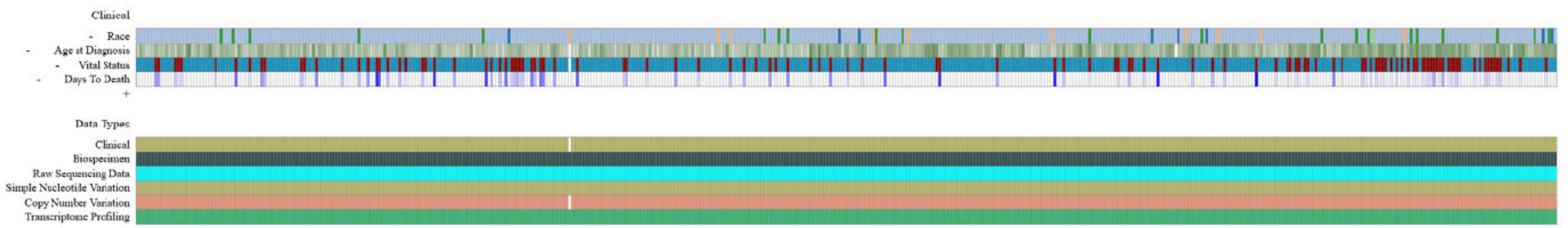
▾ *more*

**Annotation** Cancer Gene Census

| Project | Disease Type | Site | # Affected Cases |
|---|---|---|---|
| TCGA-LGG | Brain Lower Grade Glioma | Brain | 394 / 510 (77.25%) |
| TCGA-CHOL | Cholangiocarcinoma | Bile Duct | 6 / 51 (11.76%) |
| TCGA-LAML | Acute Myeloid Leukemia | Bone Marrow | 13 / 144 (9.03%) |
| TCGA-GBM | Glioblastoma Multiforme | Brain | 26 / 393 (6.62%) |
| TCGA-SKCM | Skin Cutaneous Melanoma | Skin | 22 / 469 (4.69%) |
| TCGA-UCEC | Uterine Corpus Endometrial Carcinoma | Uterus | 20 / 530 (3.77%) |
| TCGA-BLCA | Bladder Urothelial Carcinoma | Bladder | 9 / 412 (2.18%) |
| TCGA-LIHC | Liver Hepatocellular Carcinoma | Liver | 7 / 364 (1.92%) |
| TCGA-COAD | Colon Adenocarcinoma | Colorectal | 7 / 400 (1.75%) |
| TCGA-PRAD | Prostate Adenocarcinoma | Prostate | 5 / 498 (1.00%) |

# ⊪ Cancer Distribution

541 CASES AFFECTED BY 72 MUTATIONS ACROSS 24 PROJECTS

Clear | Program Name | IS | TCGA | AND | Project Id | IS | TCGA-LGG
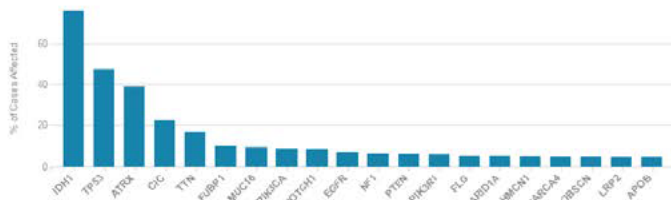
View Files in Repository

Cases (516) | Genes (14,016) | Mutations (38,973) | OncoGrid

### Genes

**Distribution of Most Frequently Mutated Genes**



Genes (x-axis): IDH1, TP53, ATRX, CIC, TTN, FUBP1, MUC16, PIK3CA, NOTCH1, EGFR, NF1, PTEN, PIK3R1, FLG, ARID1A, HMCN1, SMARCA4, OBSCN, LRP2, APOB

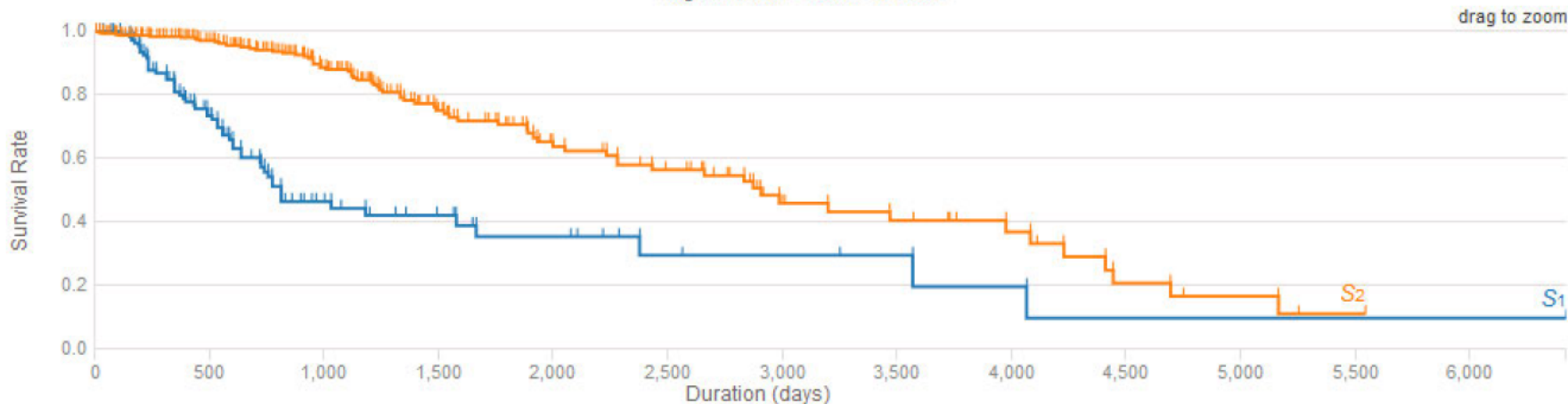*The NEW ENGLAND JOURNAL of MEDICINE*

**ORIGINAL ARTICLE**

**Overall Survival Plot**

$S_1$ (N = 122) - IDH1 Not Mutated Cases    $S_2$ (N = 389) - IDH1 Mutated Cases

Log-Rank Test P-Value = 9.20e-14



JSON | TSV | Save/Edit Gene Set

| # Affected Cases in Cohort | # Affected Cases Across the GDC | # Mutations | Annotations | Survival |
|---|---|---|---|---|
| 394 / 510 (77.25%) | 566 / 10,202 | 4 | | |

# *IDH1* and *IDH2* Mutations in Gliomas

**RESULTS**

We identified mutations that affected amino acid 132 of *IDH1* in more than 70% of WHO grade II and III astrocytomas and oligodendrogliomas and in glioblastomas that developed from these lower-grade lesions. Tumors without mutations in *IDH1* often had mutations affecting the analogous amino acid (R172) of the *IDH2* gene. Tumors with *IDH1* or *IDH2* mutations had distinctive genetic and clinical characteristics, and patients with such tumors had a better outcome than those with wild-type *IDH* genes. Each of four tested *IDH1* and *IDH2* mutations reduced the enzymatic activity of the encoded protein.

44

D-2-hydroxyglutarate

IDH1 and IDH2 mutations as novel therapeutic targets: current perspectives

# Many Types of Data Available

⊞ Summary

| | |
|---|---|
| **Project ID** | TCGA-LGG |
| **Project Name** | Brain Lower Grade Glioma |
| **Disease Type** | Brain Lower Grade Glioma |
| **Primary Site** | Brain |
| **Program** | TCGA |

Cases and File Counts by Data Category

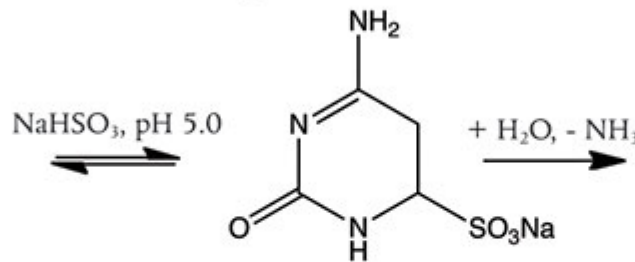| Data Category | Cases (n=516) | Files (n=12,603) |
|---|---|---|
| ■ Raw Sequencing Data | 516 | 2,105 |
| ■ Transcriptome Profiling | 516 | 2,647 |
| ■ Simple Nucleotide Variation | 513 | 4,248 |
| ■ Copy Number Variation | 514 | 2,038 |
| ■ DNA Methylation | 516 | 534 |
| ■ Clinical | 515 | 515 |
| ■ Biospecimen | 516 | 516 |

Cytosine

5-Methylcytosine (5-mC)

Step 1

**Denaturation**
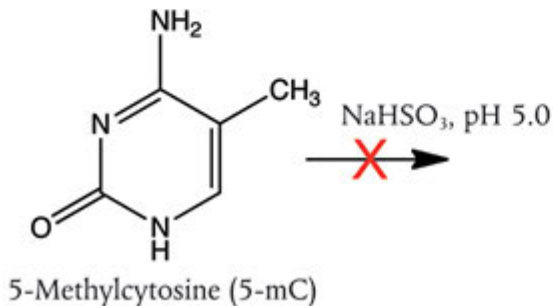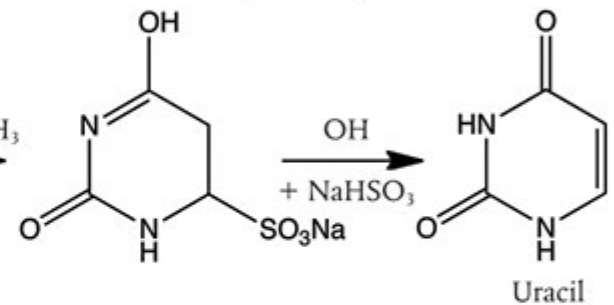Incubation at 95°C
fragments genomic DNA

Step 2

**Conversion**
Incubation with sodium bisulfite
at 65°C and low pH (5-6)
deaminates cytosine residues
in fragmented DNA

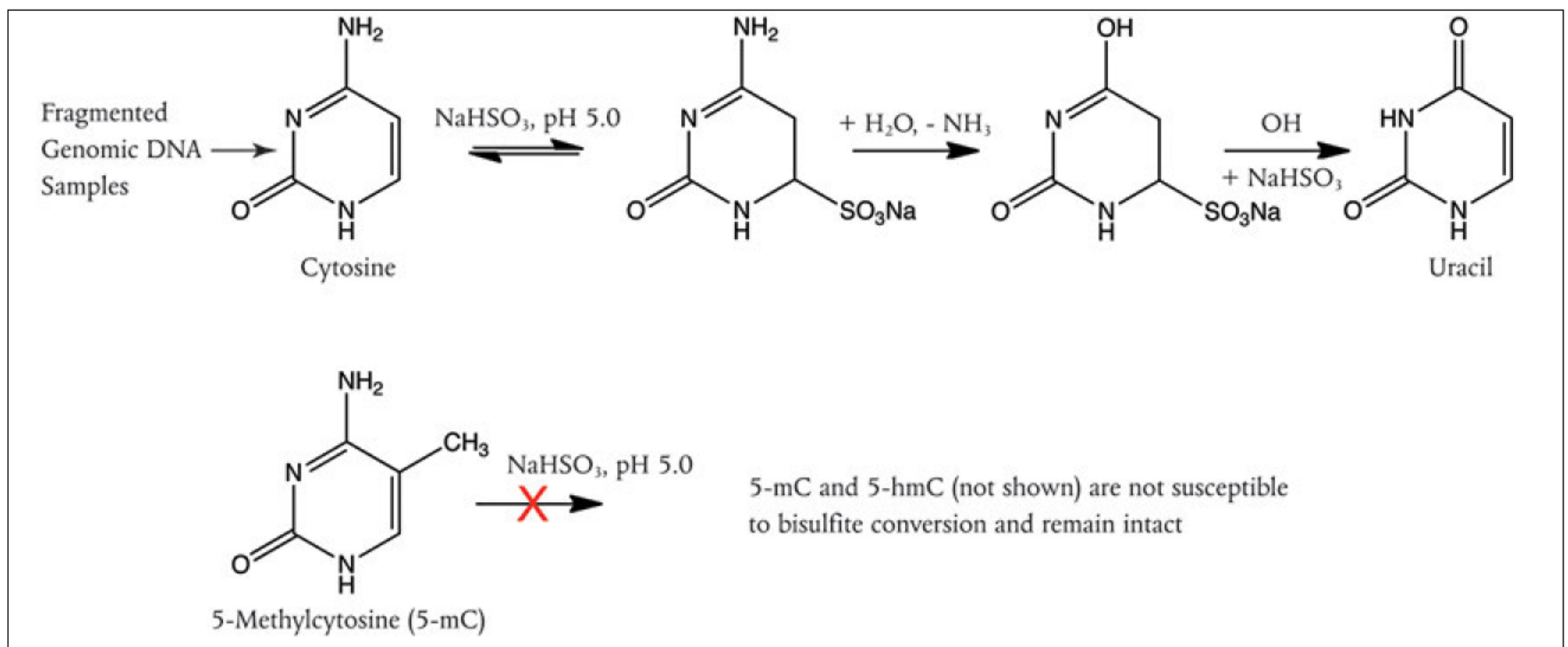Step 3

**Desulphonation**
Incubation at high pH
at room temperature for 15 min
removes the sulfite moeity,
generating uracil

Fragmented Genomic DNA Samples → Cytosine — NaHSO₃, pH 5.0 → ⇌ → (SO₃Na intermediate) — + H₂O, - NH₃ → (OH intermediate) — OH, + NaHSO₃ → Uracil

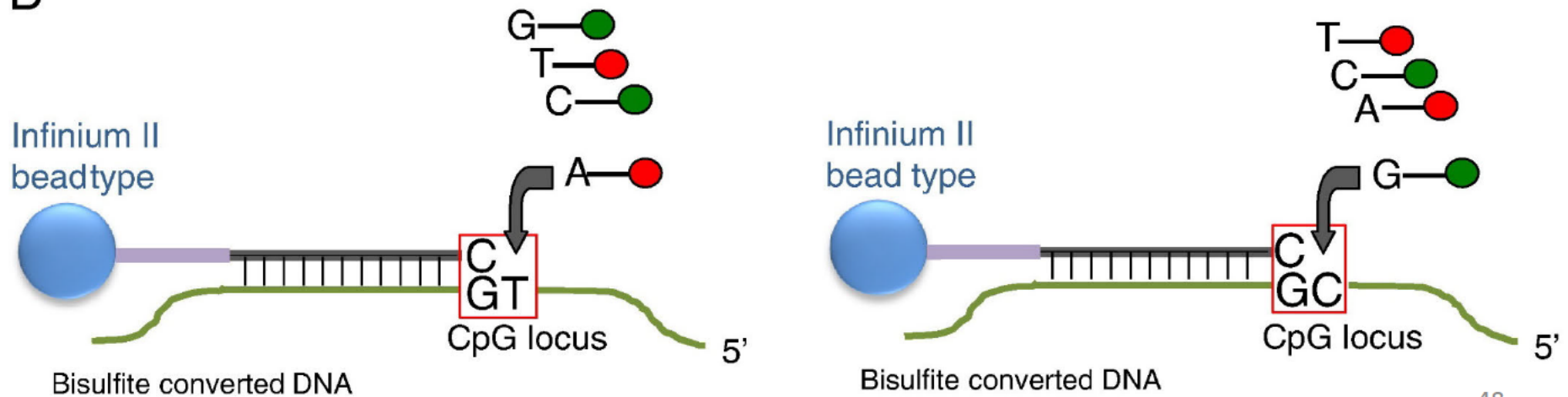5-Methylcytosine (5-mC) — NaHSO₃, pH 5.0 ✗ → 5-mC and 5-hmC (not shown) are not susceptible to bisulfite conversion and remain intact

Unmethylated locus          Methylated locus

B

Infinium II beadtype — Bisulfite converted DNA

Infinium II bead type — Bisulfite converted DNA
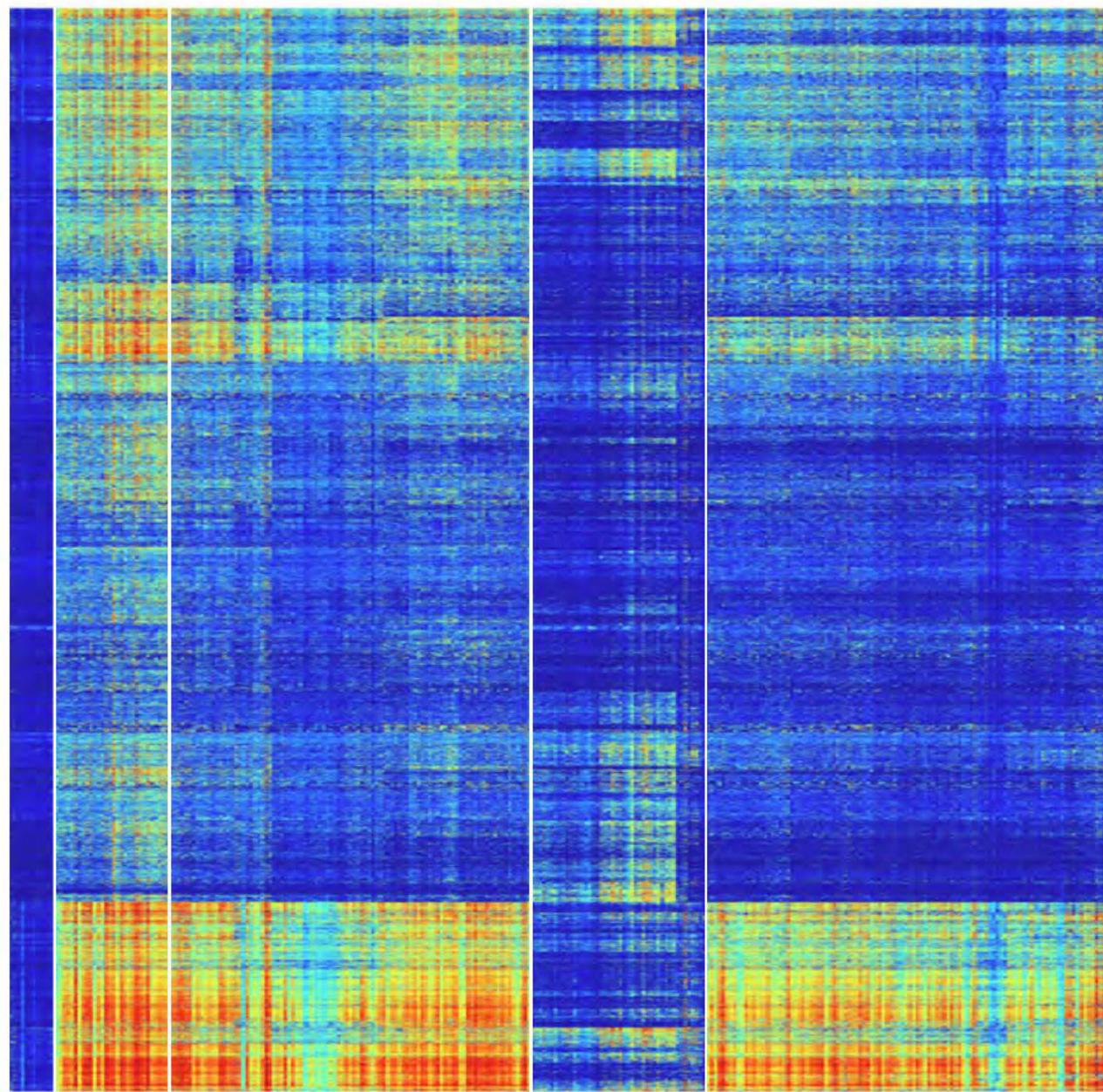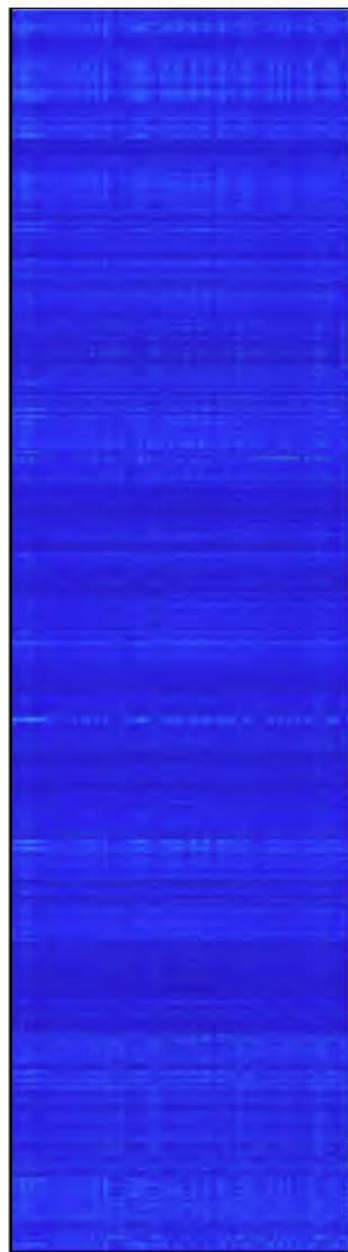
CpG locus

48

# Comprehensive, Integrative Genomic Analysis of Diffuse Lower-Grade Gliomas

The Cancer Genome Atlas Research Network*

49

Non Tumor Brain (GEO)

TCGA 289 LGG Tumors
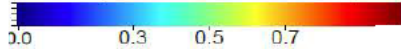
11,977 Tumor Specific CpG Probes

Methylation Beta Value for Heatmap

0.0    0.3    0.5    0.7

M1    M2    M3    M4    M5

CpG Location
CpG Island
CpG Shore

50

Groups M1 and M4 do not have IDH1 mutations. All the other groups do.

51

Groups M1 and M4 do not have IDH1 mutations. All the other groups do.

Progression Free
Kaplan-Meier Survival Curves
TCGA LGG SAMPLES (DNA Methylation Clusters)

DNA METHYLATION Clusters
— M1 (n=9)
— M2 (n=25)
— M3 (n=89)
— M4 (n=35)
— M5 (n=88)

PERCENT PROBABILITY OF SURVIVAL

TIME SINCE DIAGNOSIS (YEARS)