# TCGA

The Cancer Genome Atlas

# TCGA: History and Goal

- History:
  - Started in 2005 by the **National Cancer Institute** (NCI) and the **National Human Genome Research Institute** (NHGRI) with *$110 Million* to catalogue genetic mutations responsible for cancer (2006-2009).
  - US Government dedicated ~*$500 Million for the next 5 years* (2010-2015) to characterize **20-30 Cancers**.
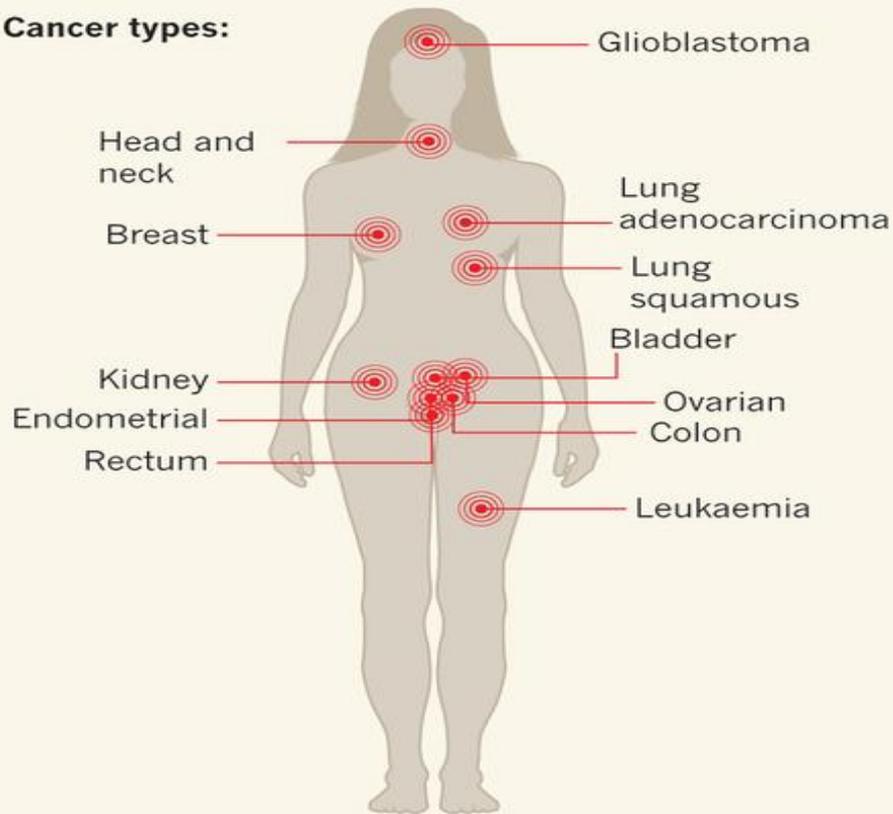
- Objective/Goal:
  - Comprehensive and coordinated effort to accelerate our **understanding** of the **molecular basis of cancer** through the application of genome analysis technologies, **including large-scale genome sequencing.**
  - To improve our ability to **diagnose, treat, and prevent cancer** through a better understanding of the **molecular basis of this disease**.

# TCGA Map - USA

**Cancer types:**



- Glioblastoma
- Head and neck
- Breast
- Lung adenocarcinoma
- Lung squamous
- Bladder
- Kidney
- Endometrial
- Rectum
- Ovarian
- Colon
- Leukaemia
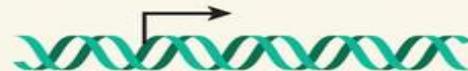
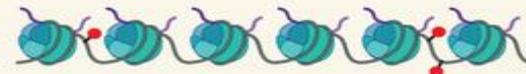**Tumour characteristics:**

DNA mutation

GATTCAT**CGT**TCCCATC

Copy-number variation

Gene expression

DNA methylation
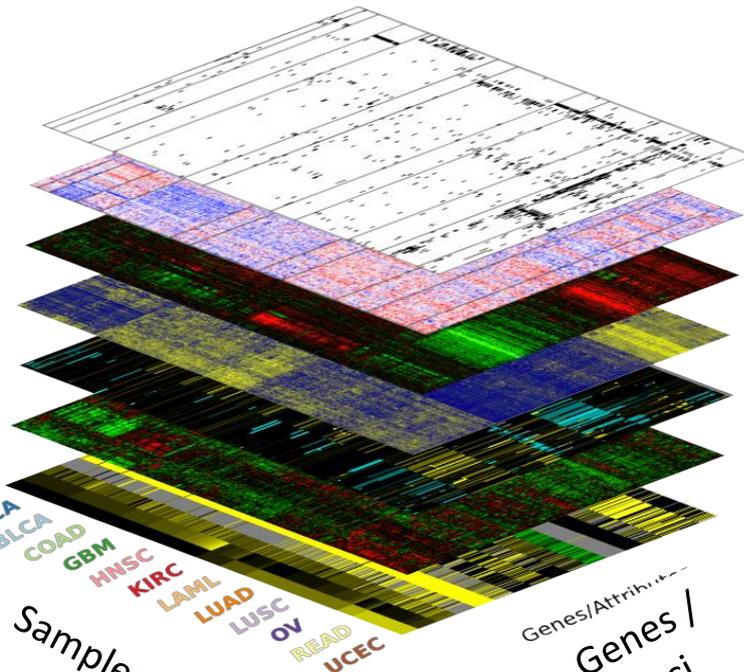
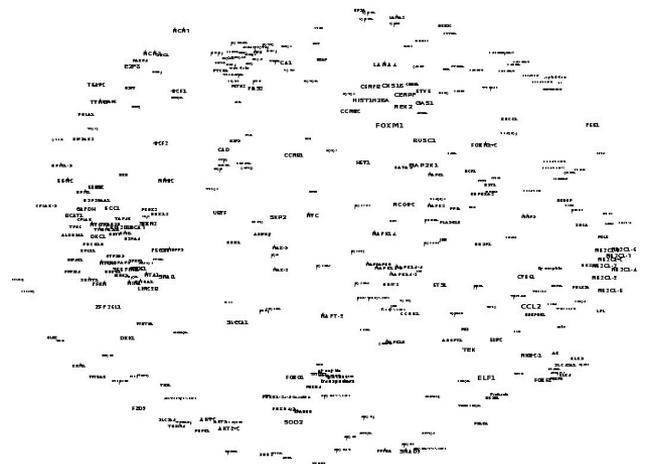MicroRNA activity

Cellular protein activity

Clinical data

Glioblastoma (GBM)
Leukemia (LAML)
Head & Neck (HNSC)
Lung Adeno (LUAD)
Lung Squamous (LUSC)
Breast (BRCA)
Kidney (KIRC)
Bladder (BLCA)
Ovarian (OV)
Endometrial (UCEC)
Colon (COAD)
Rectum (READ)

Platform

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/Attributes

Genes / Loci

# The CANCER GENOME challenge

Databases could soon be flooded with genome sequences from 25,000 tumours. **Heidi Ledford** looks at the obstacles researchers face as they search for meaning in the data.
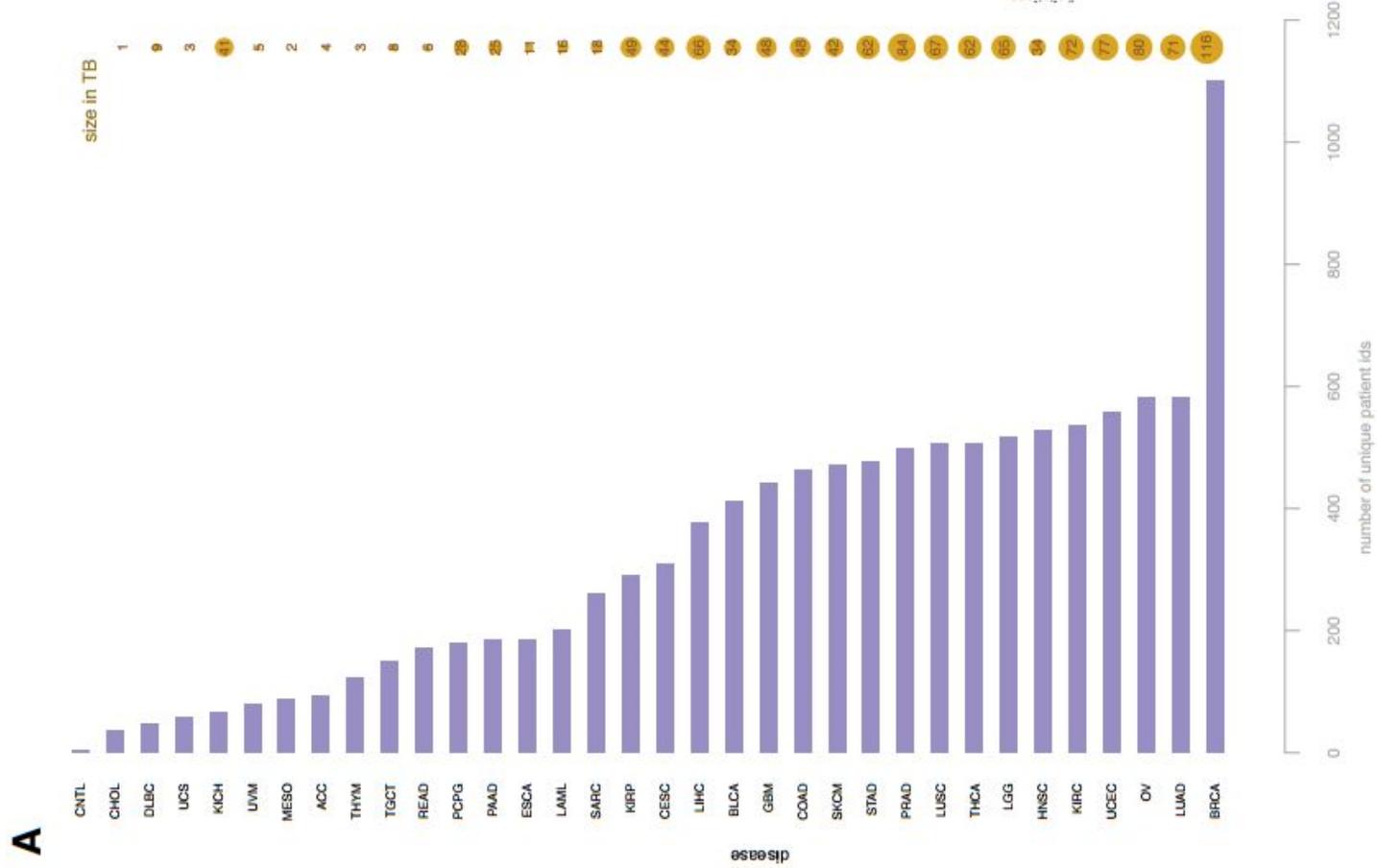
W hen it was first discovered, in 2006, in a study of 35 colorectal cancers[1], the mutation in the
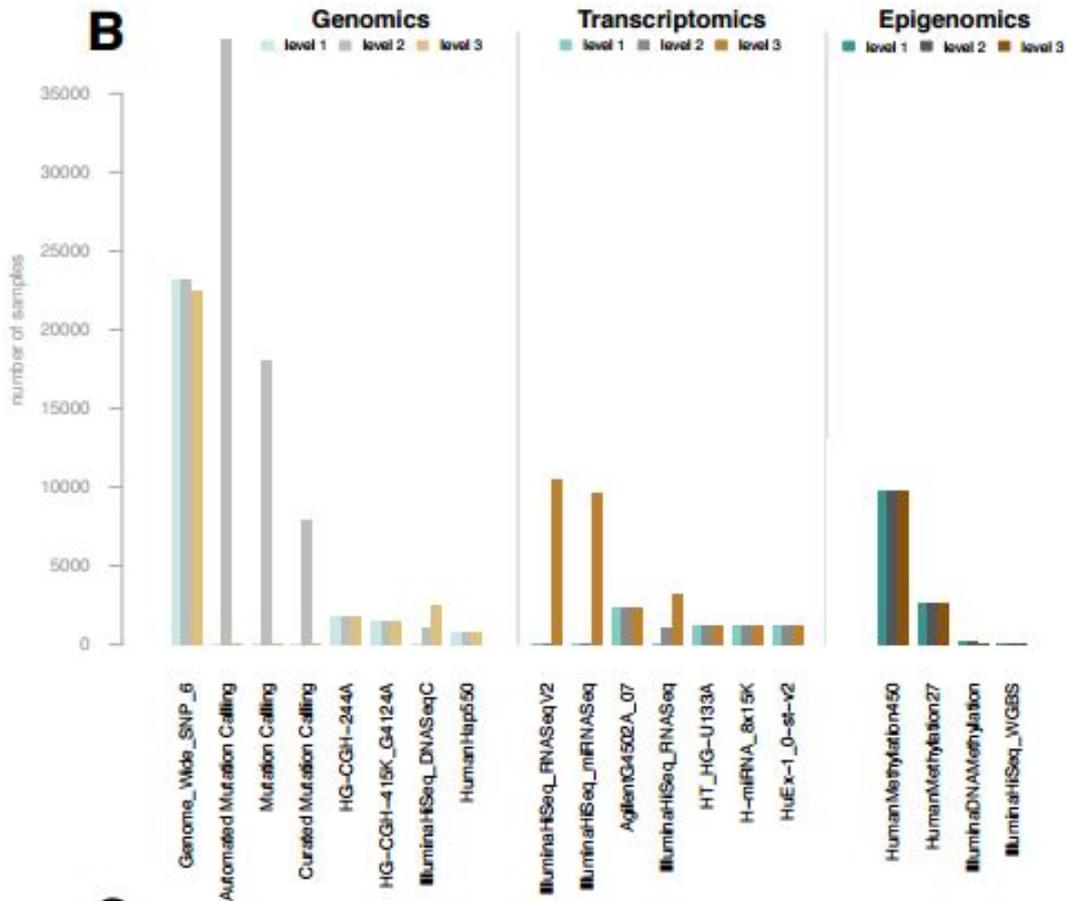
## GENOMES AT A GLANCE

Circos plots can give a snapshot of the mutations within

needle pulled from a veritable haystack of cancer-associated mutations thanks to high-powered genome sequencing. In the past two

# Number of TCGA Samples

# TCGA Molecular Data

# TCGA Impact



# publications (TCGA Research Network)

18 (1)    21 (0)    38 (2)    62 (1)    127 (7)    252 (8)    413 (7)    320 (3)

# citations (left axis): 8000, 4000, 0

# citations TCGA / # citations cancer (right axis): 0.003, 0.002, 0.001, 0.000

Bar values: 12 (2008), 294 (2009), 690 (2010), 1152 (2011), 2384 (2012), 4582 (2013), 7256 (2014), 5116 (2015*)

Years: 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015*

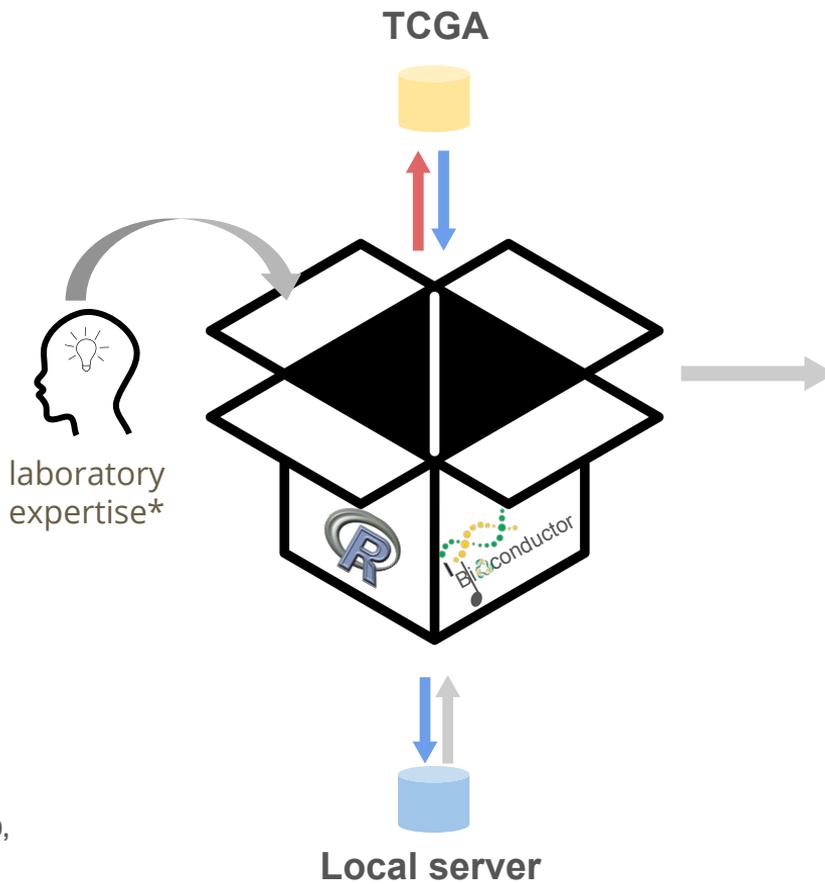* until August 2015

# Problems

i) Difficulty finding the desired information

     i.i) If we want an old version of the data the only solution is via the http

ii) How the data is organized in the portal

iii) Download complexity of the site

# Solution

TCGA

Search
Retrieve
Analysis
Local Database
Public Database

laboratory expertise*

Local server

* PMID: 24885402, 24120142, 23717510,
22684628, 22479200, 22187159,
22120008,
21659424, 20399149

# TCGAbiolinks

An R/Bioconductor package for integrative analysis with TCGA data

# Preparing for tutorial

- /dados/uruguay/aluno/tcgabiolinks
- /dados/uruguay/aluno/elmer

# Aim

i) facilitate the TCGA open-access data retrieval

ii) prepare the data using the appropriate pre-processing strategies

iii) provide the means to carry out different standard analyses

iv) allow the user to download a specific version of the data and thus to easily reproduce earlier research results.

# Pipeline

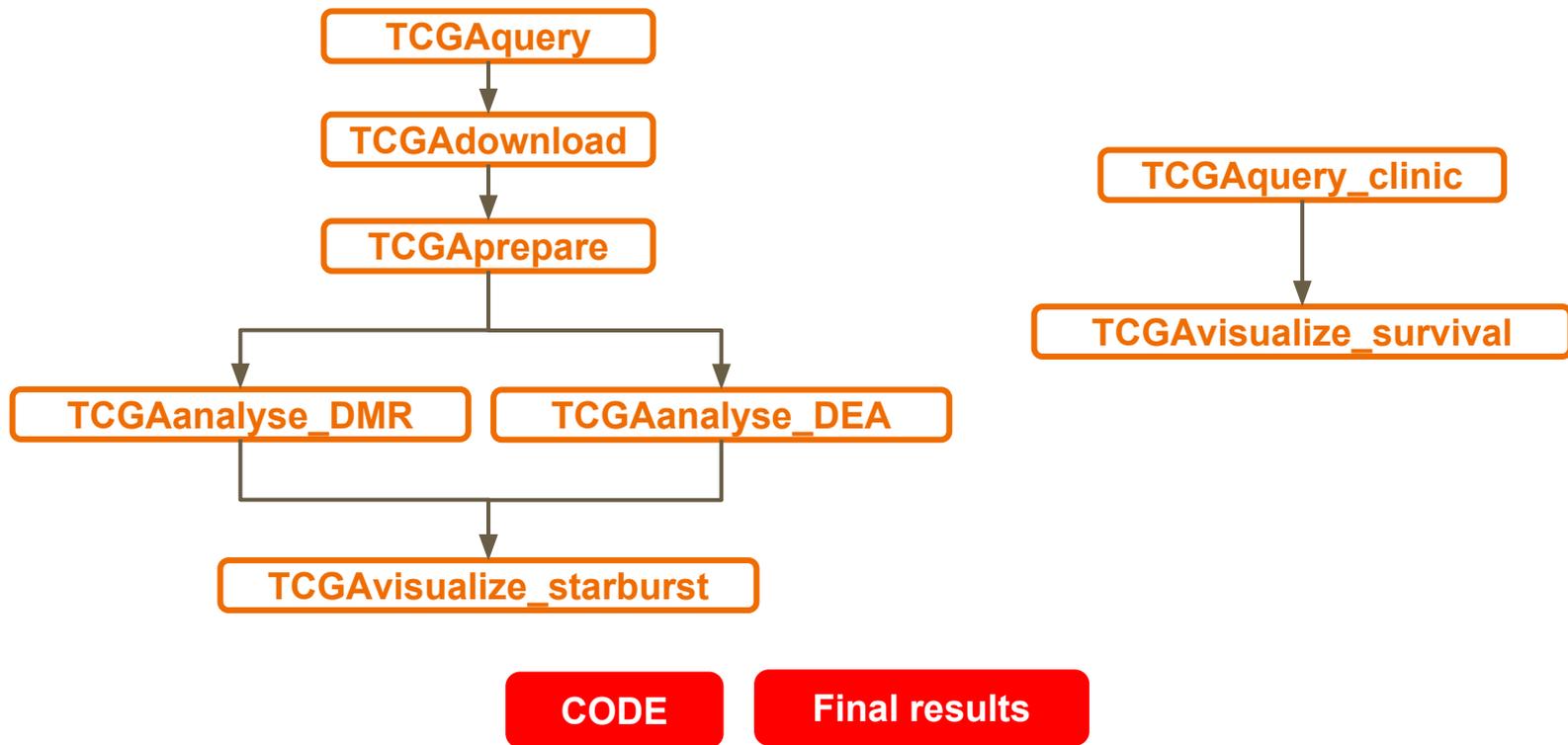# Searching for data



```
query <- TCGAquery ( tumor = "GBM",                                    # Vector of tumors
                     platform = "Humanmethylation450",                 # Vector of platforms
                     level = 3,                                        # 1, 2, 3
                     version = list(c("HumanMethylation450","GBM",5)),  # List of triple (tumor, platform, version)
                     samples = c("TCGA-06-6694-01A-12D-1844-05",       # Vector of barcodes
                                 "TCGA-06-0171-02A-11D-2004-05")
)
```

# Searching for data - Example

query <- TCGAquery ( tumor = "GBM", platform = "Humanmethylation450", level = 3)

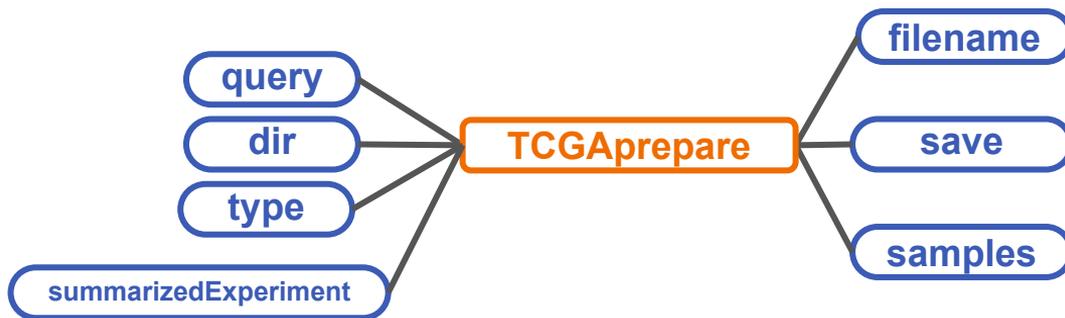| addedDate | baseName | barcode | name | revision |
|---|---|---|---|---|
| 2014-11-04 | jhu-usc.edu_GBM_HumanMethylation450 | TCGA-06-6694-01A-12D-1844-05,TCGA-74-6578-01... | jhu-usc.edu_GBM.HumanMethylation450.Level_3.8.... | 6 |
| 2014-11-04 | jhu-usc.edu_GBM_HumanMethylation450 | TCGA-06-0171-02A-11D-2004-05,TCGA-06-0190-02... | jhu-usc.edu_GBM.HumanMethylation450.Level_3.2.... | 6 |
| 2014-11-04 | jhu-usc.edu_GBM_HumanMethylation450 | TCGA-19-4065-02A-11D-2004-05,TCGA-19-4065-01... | jhu-usc.edu_GBM.HumanMethylation450.Level_3.9.... | 6 |
| 2014-11-04 | jhu-usc.edu_GBM_HumanMethylation450 | TCGA-06-0125-02A-11D-2004-05,TCGA-06-0125-01... | jhu-usc.edu_GBM.HumanMethylation450.Level_3.1.... | 6 |
| 2014-11-04 | jhu-usc.edu_GBM_HumanMethylation450 | TCGA-76-4932-01A-01D-1481-05,TCGA-28-5214-01... | jhu-usc.edu_GBM.HumanMethylation450.Level_3.6.... | 6 |
| 2014-11-04 | jhu-usc.edu_GBM_HumanMethylation450 | TCGA-07-0227-20A-01D-A368-05,TCGA-06-AABW-1... | jhu-usc.edu_GBM.HumanMethylation450.Level_3.11... | 6 |
| 2014-11-04 | jhu-usc.edu_GBM_HumanMethylation450 | TCGA-19-0957-02A-11D-2004-05,TCGA-19-1389-02... | jhu-usc.edu_GBM.HumanMethylation450.Level_3.5.... | 6 |
| 2014-11-04 | jhu-usc.edu_GBM_HumanMethylation450 | TCGA-87-5896-01A-01D-1697-05,TCGA-76-6191-01... | jhu-usc.edu_GBM.HumanMethylation450.Level_3.7.... | 6 |
| 2014-11-04 | jhu-usc.edu_GBM_HumanMethylation450 | TCGA-14-1034-02B-01D-2004-05,TCGA-14-1402-02... | jhu-usc.edu_GBM.HumanMethylation450.Level_3.4.... | 6 |
| 2014-11-04 | jhu-usc.edu_GBM_HumanMethylation450 | TCGA-06-0152-02A-01D-2004-05,TCGA-06-0152-01... | jhu-usc.edu_GBM.HumanMethylation450.Level_3.3.... | 6 |
| 2014-11-04 | jhu-usc.edu_GBM_HumanMethylation450 | TCGA-07-0227-20A-01D-A33U-05,TCGA-OX-A56R-0... | jhu-usc.edu_GBM.HumanMethylation450.Level_3.10... | 6 |
| 2014-11-04 | jhu-usc.edu_GBM_HumanMethylation450 | TCGA-07-0227-20A-01D-A392-05,TCGA-26-A7UX-01... | jhu-usc.edu_GBM.HumanMethylation450.Level_3.12... | 6 |

# Downloading the data



```
TCGAdownload(data = query,                                    # TCGAquery result
             path = ".",                                      # Path to save files
             samples = c("TCGA-06-6694-01A-12D-1844-05",      # Vector of barcodes to download file
                         "TCGA-06-0171-02A-11D-2004-05"),
             force = FALSE)                                    # If already downloaded download it again?
```

# Reading the data



```
TCGAprepare(query = query,                                # TCGAquery result
            dir = ".",                                    # Path where files were saved
            samples = c("TCGA-06-6694-01A-12D-1844-05",   # Vector of barcodes to read file
                        "TCGA-06-0171-02A-11D-2004-05"),
            save = TRUE,                                  # Save prepared object?
            filename = "name.rda",                        # Name of the file with the prepared object
            summarizedExperiment = TRUE )                 # If FALSE output is a data.frame
```

# SummarizedExperiment



Huber, Wolfgang, et al. "Orchestrating high-throughput genomic analysis with Bioconductor." *Nature methods* 12.2 (2015): 115-121.

# Results



**Mean DNA methylation**

hypermutated  ● 0  ▲ 1  ■ NA

**Groups** ■ CIMP.H (n = 27) ■ CIMP.L (n = 26)

# Results



**GO:Molecular Function**

0.0

- olfactory receptor activity (n=6)
- transition metal ion binding (n=856)
- zinc ion binding (n=832)
- purine nucleotide binding (n=765)
- nucleotide binding (n=872)
- ribonucleotide binding (n=739)
- purine ribonucleotide binding (n=739)
- ion binding (n=923)
- cation binding (n=920)
- metal ion binding (n=920)

−log10(FDR)

**Pathways**

0.0

- Molecular Mechanisms of Cancer (n=318)
- Protein Ubiquitination Pathway (n=238)
- EIF2 Signaling (n=160)
- Role of Macrophages, Fibroblasts and Endothelial Cells in Rheumatoid Arthritis (n=255)
- mTOR Signaling (n=171)
- Regulation of eIF4 and p70S6K Signaling (n=134)
- Glucocorticoid Receptor Signaling (n=234)
- Mitochondrial Dysfunction (n=135)
- B Cell Receptor Signaling (n=149)
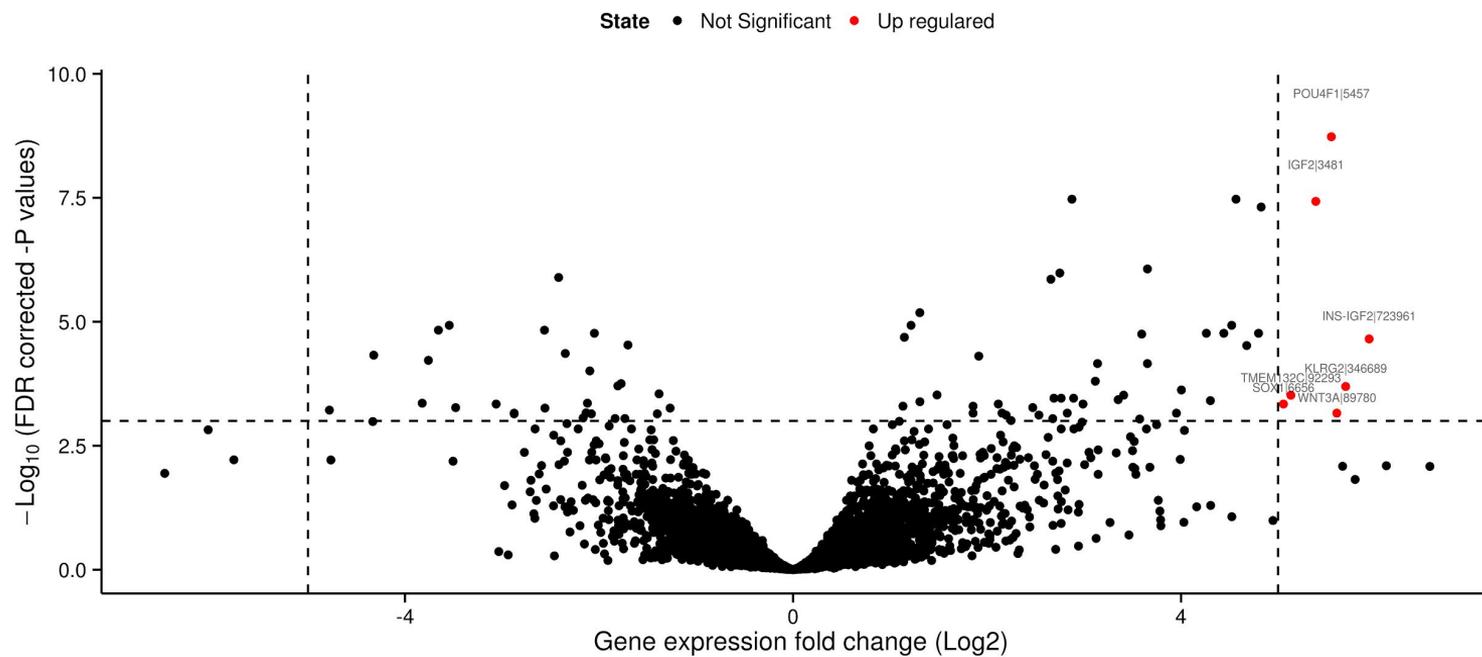- Superpathway of Inositol Phosphate Compounds (n=168)

−log10(FDR)

# Results


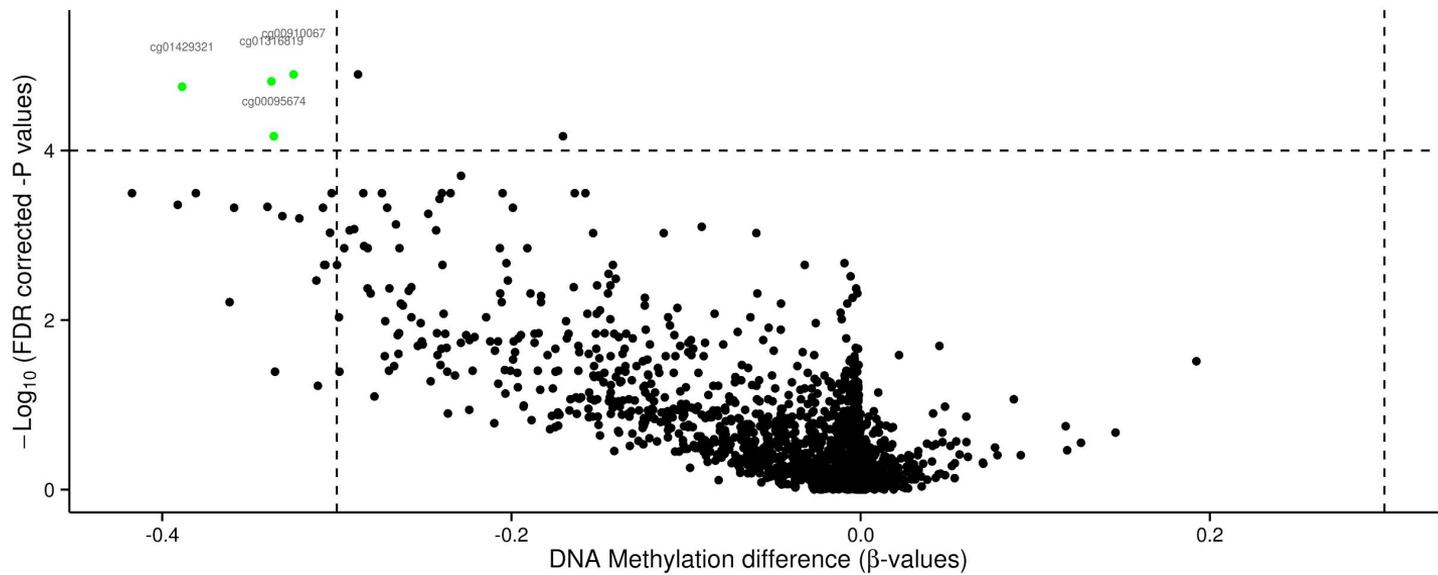
Volcano plot (CIMP.I vs CIMP.H)
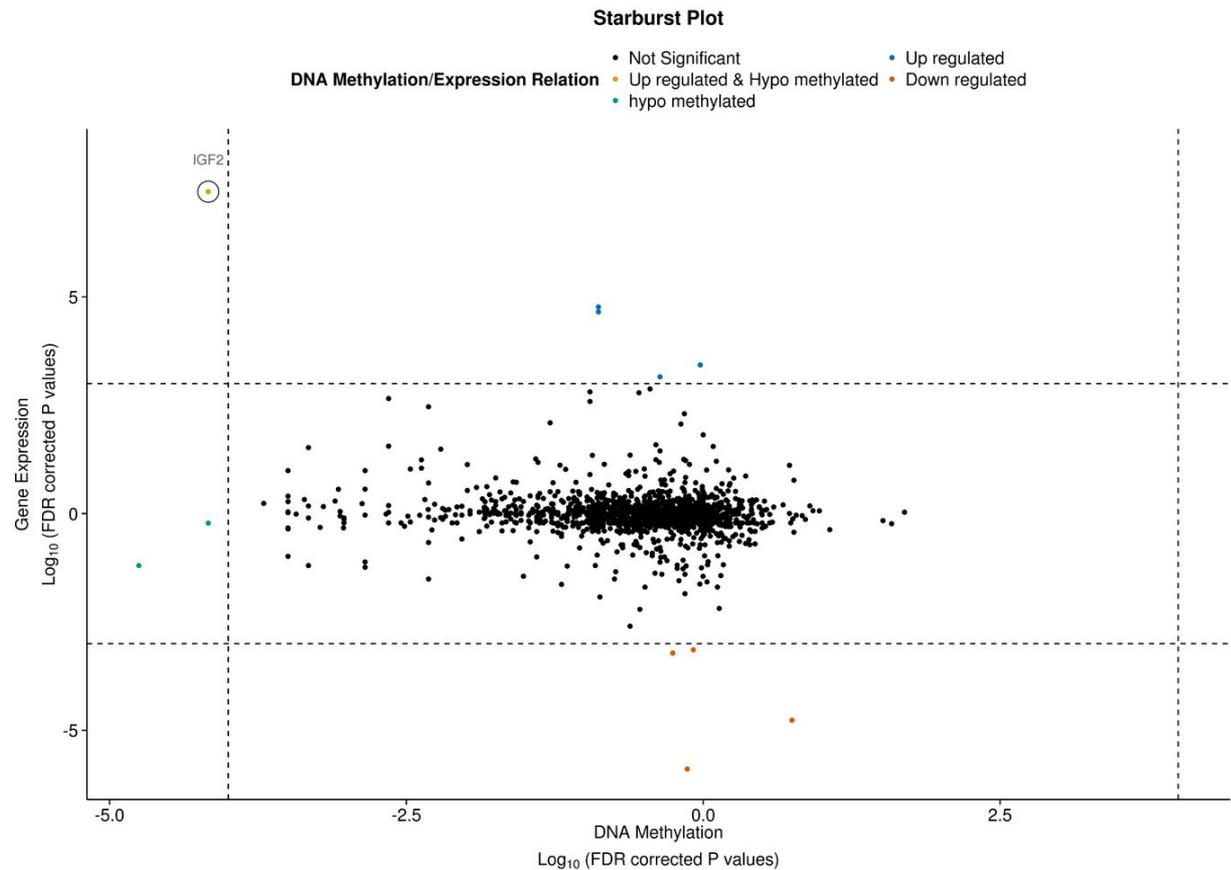
# Results



Volcano plot ( CIMP.L vs CIMP.H )

State ● Not Significant ● Hypomethylated in CIMP.L

# Starburst



**Starburst Plot**

**DNA Methylation/Expression Relation**
- Not Significant
- Up regulated & Hypo methylated
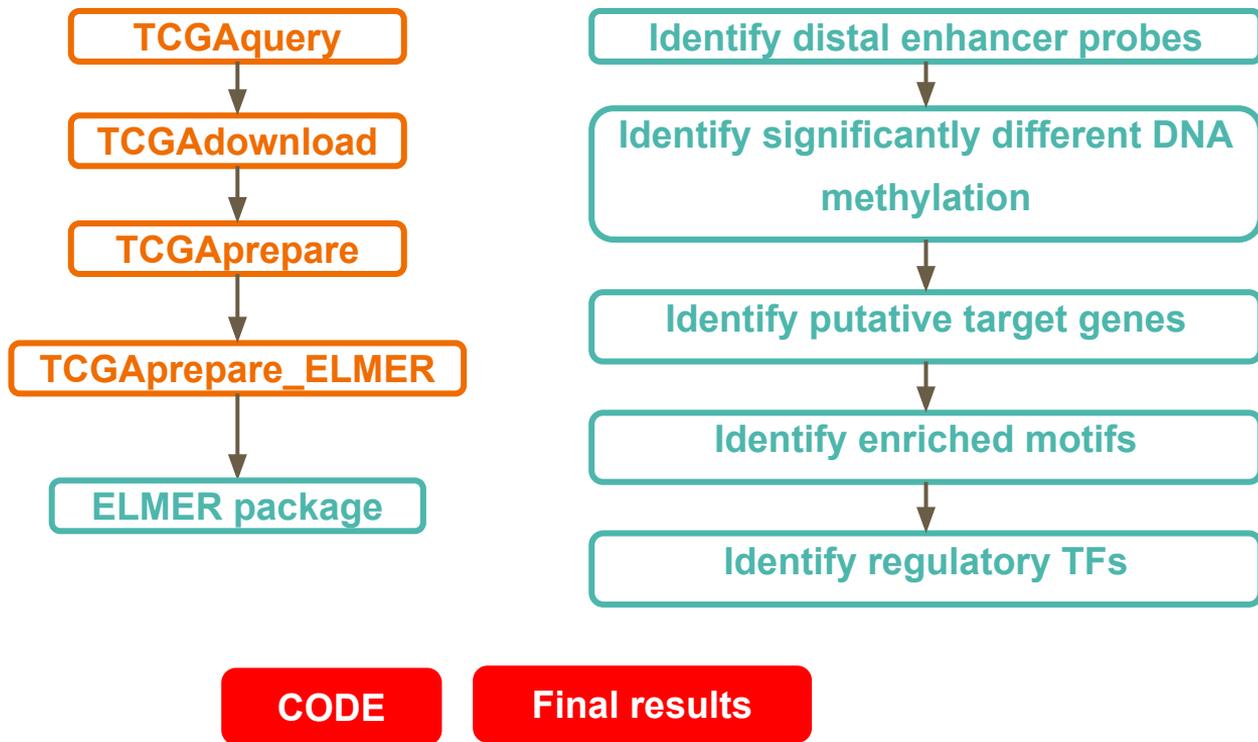- hypo methylated
- Up regulated
- Down regulated
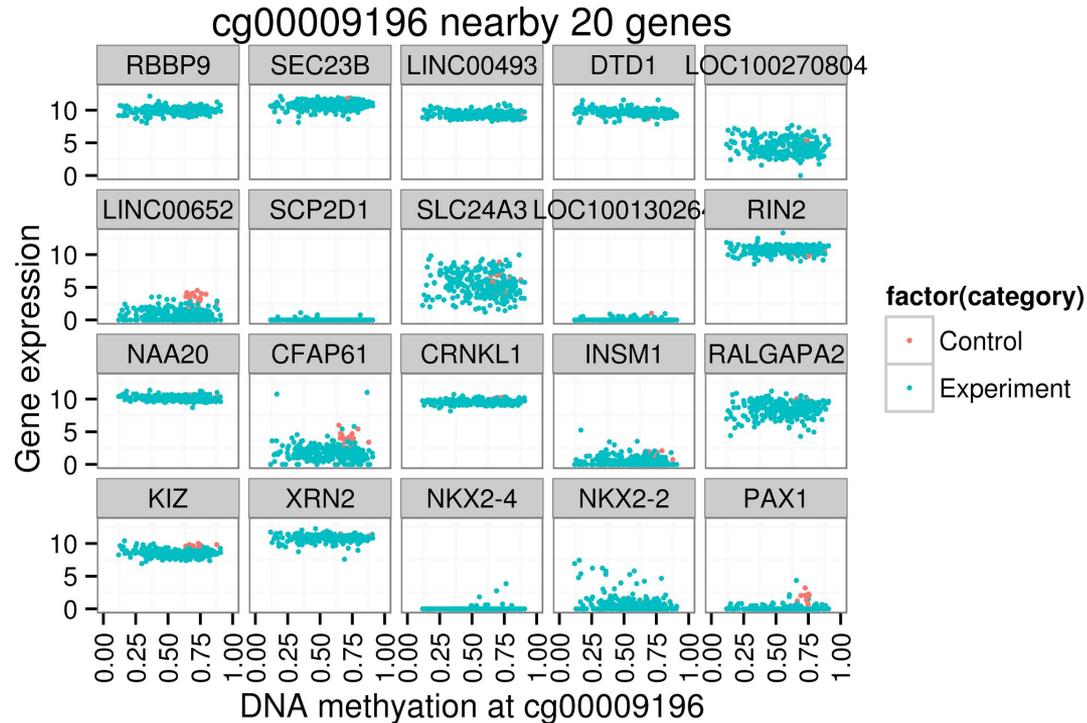
# Results

# ELMER package

ELMER is designed to use DNA methylation and gene expression from a large number of samples to infere regulatory element landscape and transcription factor network in primary tissue.

# Pipeline - TCGAbiolinks + ELMER

```
TCGAquery
   ↓
TCGAdownload
   ↓
TCGAprepare
   ↓
TCGAprepare_ELMER
   ↓
ELMER package
```

```
Identify distal enhancer probes
   ↓
Identify significantly different DNA methylation
   ↓
Identify putative target genes
   ↓
Identify enriched motifs
   ↓
Identify regulatory TFs
```

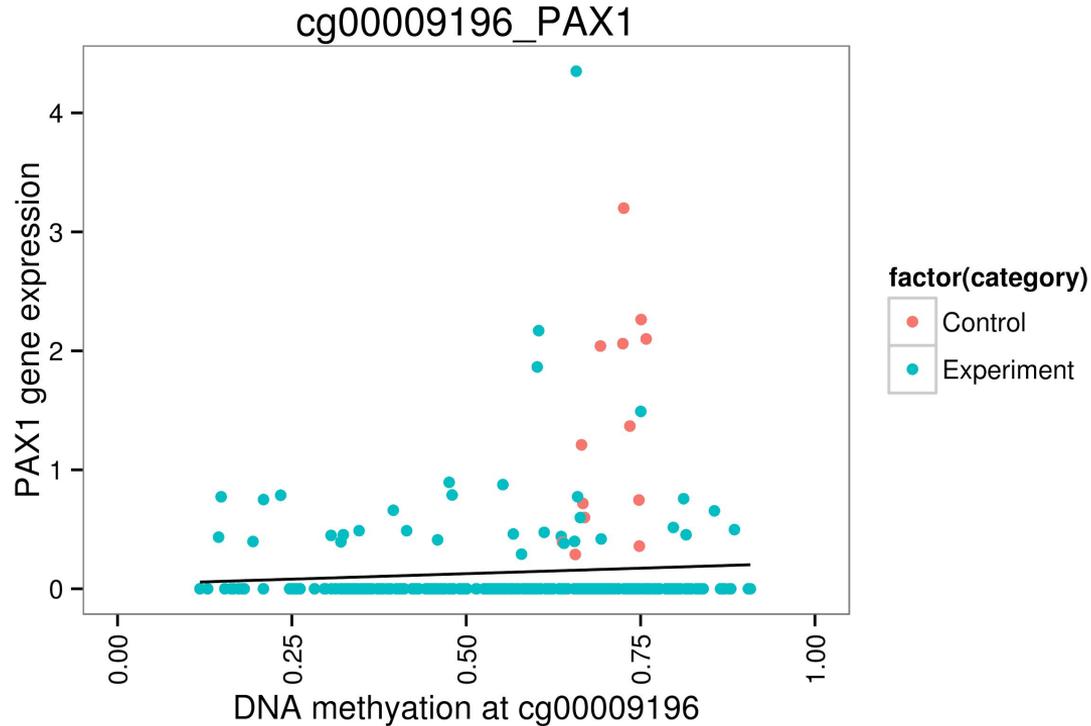**CODE**   **Final results**

# Results

Generate scatter plots for one probes' nearby 20 gene expression vs DNA methylation at this probe.

# Results

You can also focus on one probe-gene pair. The entrez gene ID for PAIX is 5075.
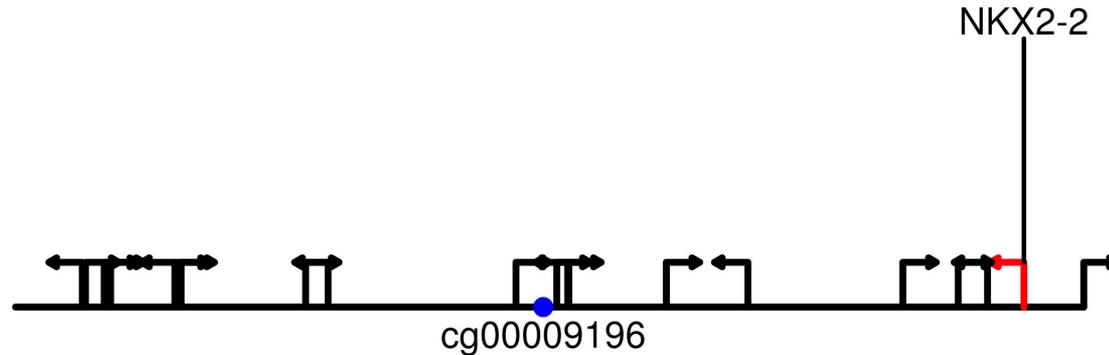


cg00009196_PAX1

# Results

You can generate schematic plot for one probe with 20 nearby genes and label the gene significantly linked with the probe in red.
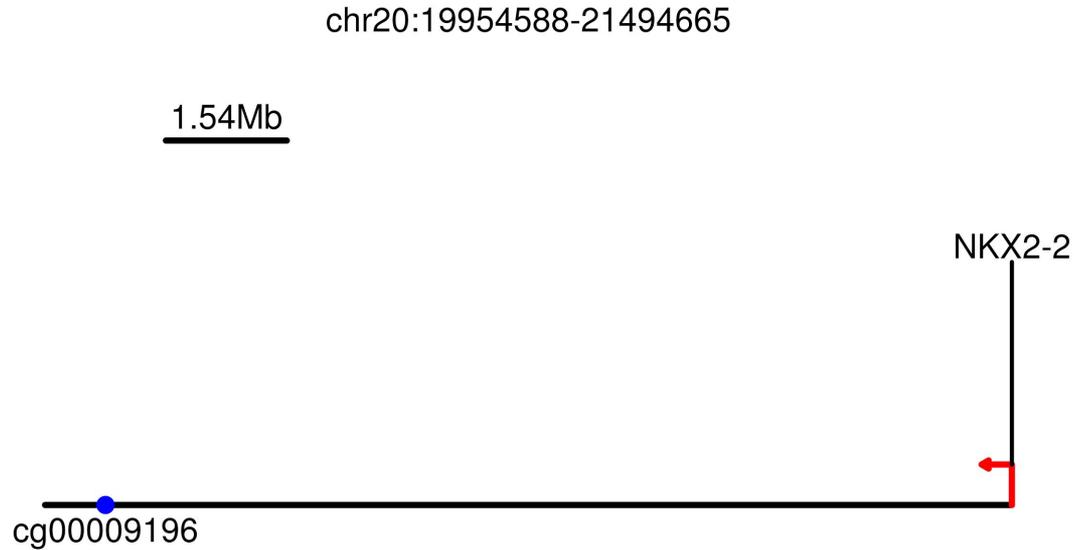
chr20:18477888-21686297

3.21Mb

NKX2-2

cg00009196

# Results

Generate schematic plot for one gene with the probes which the gene is significantly linked to.

chr20:19954588-21494665

1.54Mb

NKX2-2

cg00009196

# Results

Sometimes there is more than one probe nearby.

chr5:149980642-150169781

0.19Mb

SYNPO

cg04604050

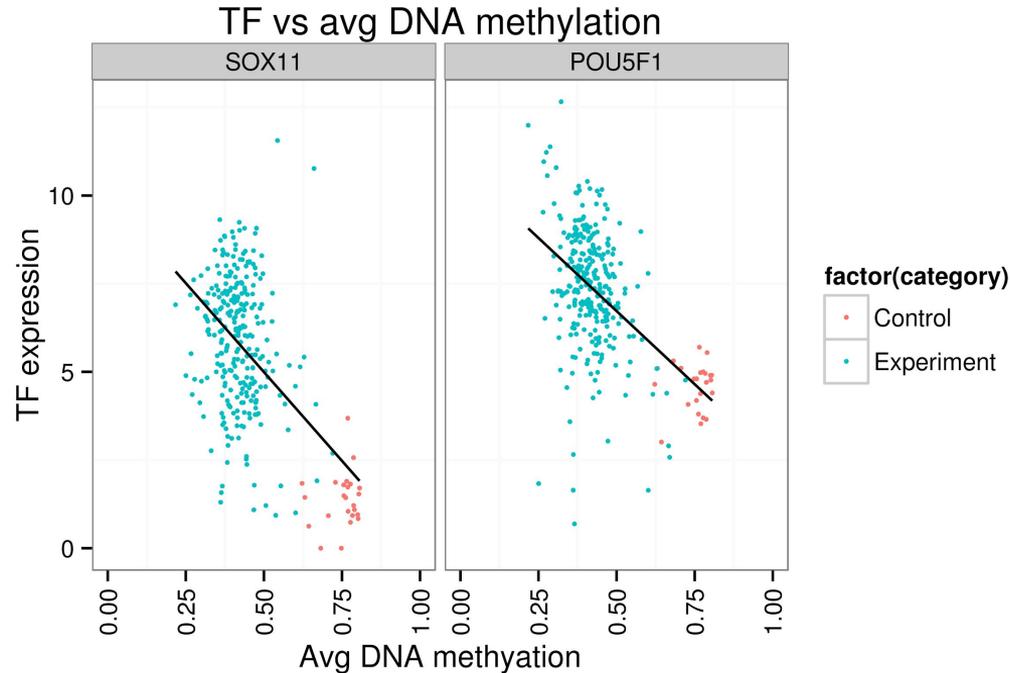cg09501687

# Results

We can generate a scatter plot for TF expression vs average DNA methylation of the sites with certain motif.

# References

- TCGAbiolinks
- ELMER

# Troubleshoot

- TCGAdownload is not working or it crashed.
    - TCGA data portal might be down or the access limit was reached and your IP was blocked. Try again in some minutes
    - If it is not working the following file has all the objects you need for the course, just upload this to your rstudio.