

TCGA  
GBM DNA Methylation Update  
Patient Centric Analysis

Houtan Noushmehr

Peter W. Laird

Jan 13, 2012

# Rules for DNA methylation Calls/Scores

- 1<sup>st</sup> step, collapse multiple CpG to 1 gene for each platform (27k, 450k).
- Using Gene Expression Data (1 gene = 1 Exp. Value per sample)
- Merged Samples & Probes (DNA methylation) with Gene Exp for each platform
  - Cataloged probes to gene by annotation provided by Illumina & collected CpGs within 500bp of a known TSS.
- Calculated spearman correlation for all CpG probes to 1 gene expression.
- Among multiple CpG, selected 1 CpG with the lowest correlation value.
- Reduce platforms from N:1, to 1:1.

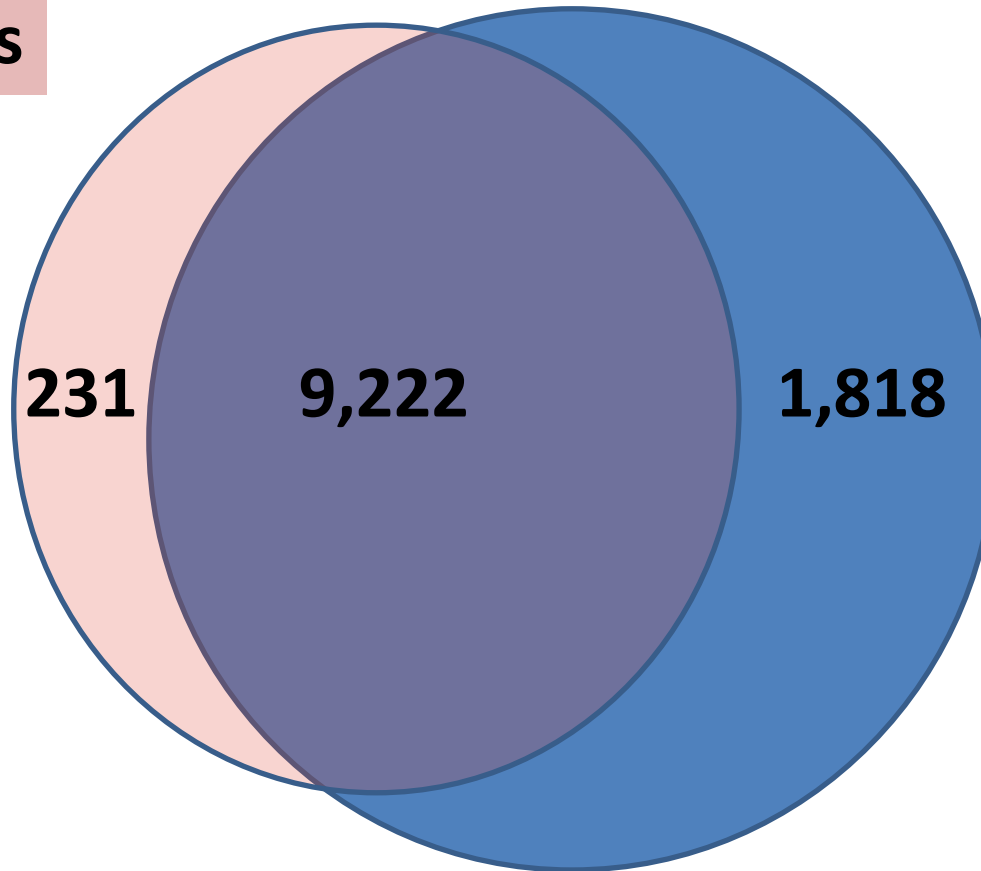
# Venn diagram DNA meth vs Gene Exp.

- TODO: add venn diagram showing the number of probes/genes overlapping the two different experiments and the different platforms.

# Genes Overlap 27K and 450K Platform

**27K**  
**9,453 Genes**

**450K**  
**11,040 Genes**



# Data available for analysis

- 27K: 279 samp X 9,453 CpG:Gene
- 450K: 74 samp X 11,040 CpG:Gene

# Rules for DNA methylation Calls/Scores

- 2<sup>nd</sup> step, create cut-offs.
- Using spearman correlation rho (GeneExp vs DNA methylation) we labeled each gene as:

SNC (Strongly Negatively Correlated)

$< -0.5$

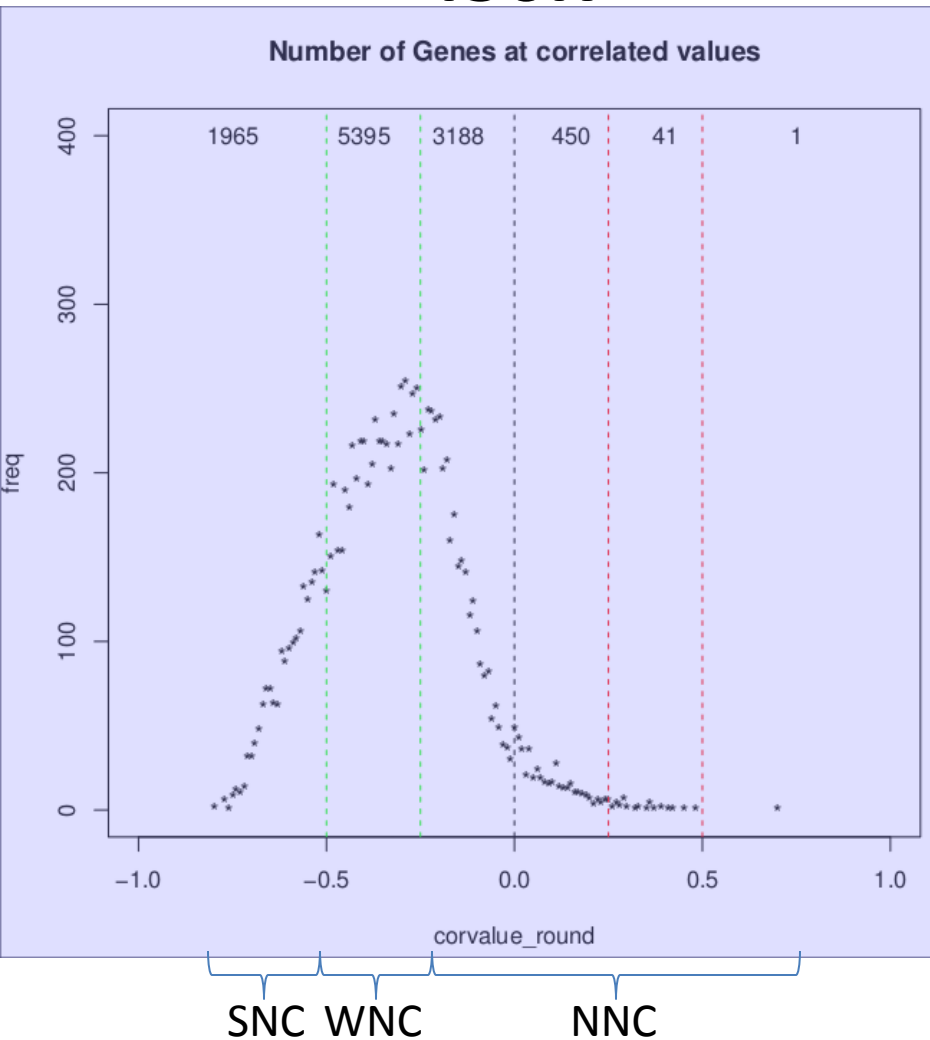
WNC (Weakly Negatively Correlated)

$-0.5$  to  $-0.25$

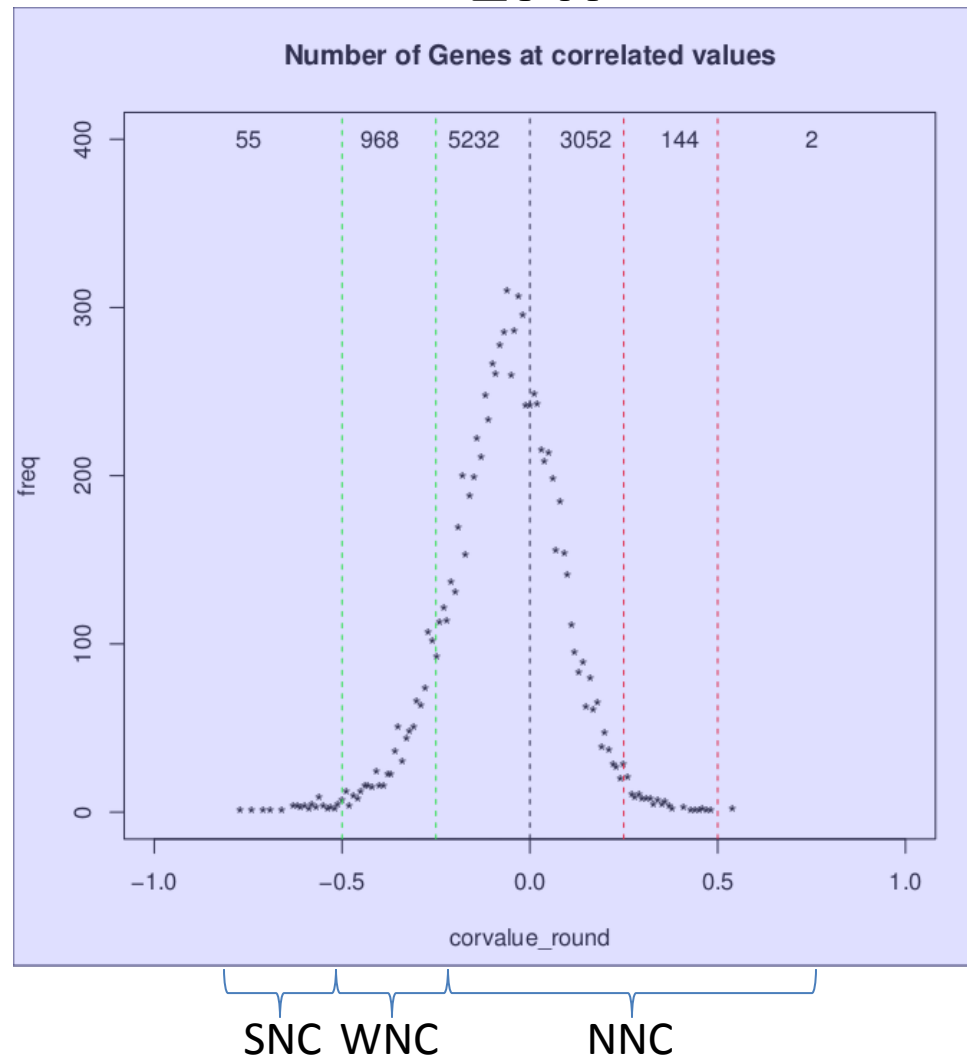
NNC (No Negative Correlation)

$> -0.25$

# 450K



# 27K



# Rules for DNA methylation Calls/Scores

- 3<sup>rd</sup> step,
  - We took the 10<sup>th</sup> and 90<sup>th</sup> percentile beta value across tumor samples for each gene per platform and labeled it:
  - **T10, T90.**
  - We did the same for normal tissues
  - **N10, N90.**
    - When analyzing the 27K platform, we used 4 non-tumor brains
    - When analyzing the 450K platform, we used 24 samples across 3 different normal tissues (72 samples):  
breast, kidney and lung
- Specifically for the 450K, we calculated the 10<sup>th</sup>, 90<sup>th</sup> for each tissue type, then calculated the median across gene for each category: 10<sup>th</sup>, 90<sup>th</sup>. This is done, so there is no influence on tissue type over another.



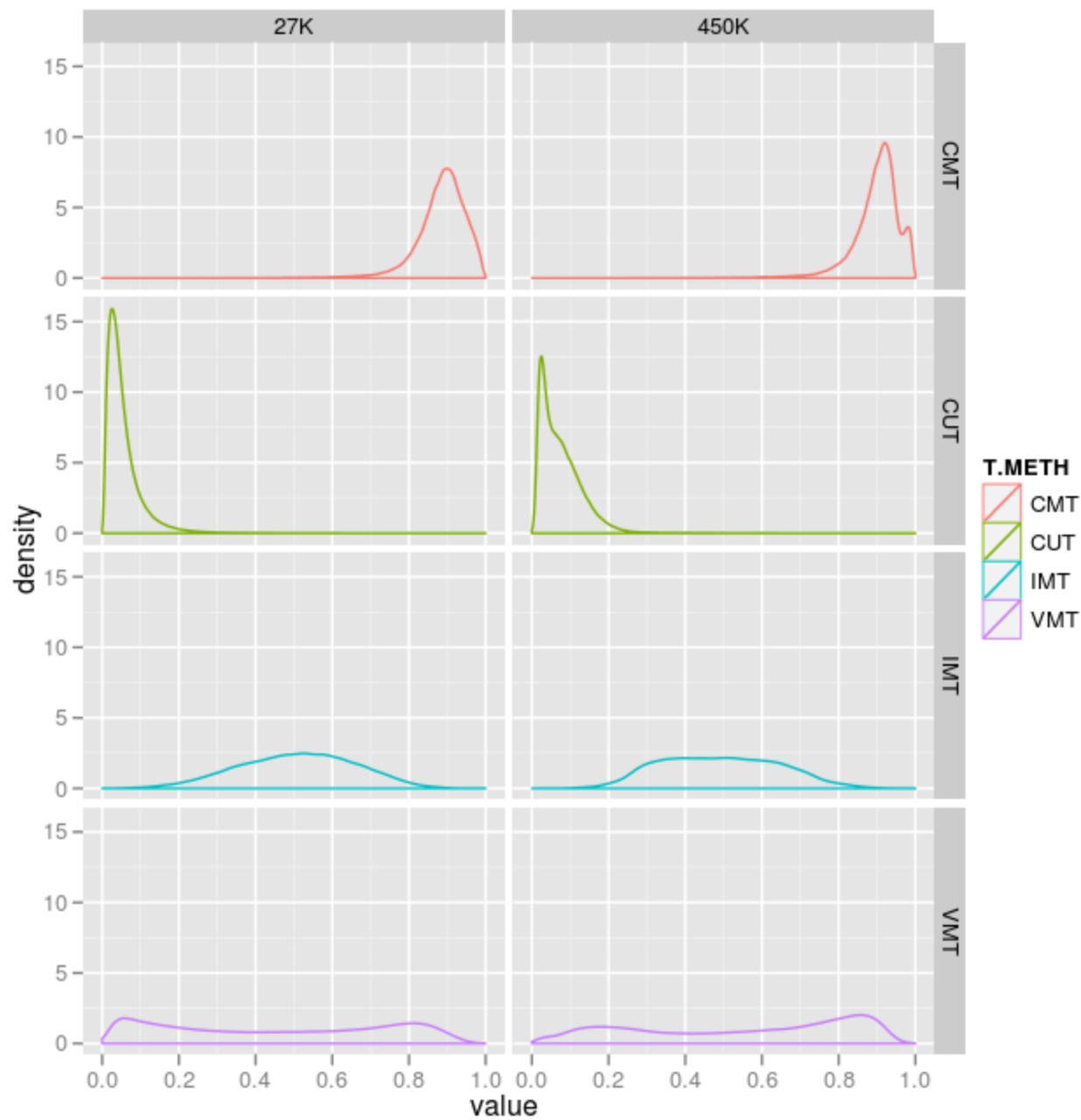
# Normal Samples

	Tumor	MatchedNormal127	UnmatchedNormal127	MatchedNormal1450	UnmatchedNormal1450
1	AML	0	0	0	0
2	BLCA	0	0	11	0
3	LGG	0	0	2	0
4	BRCA	27	0	90	0
5	CESC	0	0	0	0
6	COAD	37	3	0	0
7	GBM	0	4	0	0
8	HNSC	0	0	0	0
9	KIRC	199	0	147	5
10	KIRP	5	0	0	0
11	LIHC	0	0	0	0
12	LUAD	24	0	0	0
13	LUSC	27	0	40	2
14	OV	4	8	0	0
15	PAAU	0	0	0	0
16	PRAD	0	0	0	0
17	READ	5	2	0	0
18	STAD	43	0	2	0
19	THCA	0	0	0	0
20	UCEC	1	1	15	11

Randomly  
Selected 24 normals  
from each tumor pair  
(BRCA, KIRC, LUSC)

# Rules for DNA methylation Calls/Scores

- 4<sup>th</sup>,
- Using the 10<sup>th</sup>, 90<sup>th</sup> cut-offs for tumor and normals we defined each gene per platform per tissue type (normal and tumor) according to the following rules:
  - If percentile 90 < 0.25 → **CUN or CUT** (constitutively unmethylated in normal or tumor)
  - If percentile 10 > 0.75 → **CMN or CMT** (constitutively methylated in normal or tumor)
  - If percentile 10 > 0.25 && percentile 90 < 0.75 → **IMN or IMT** (intermediate methylated in normal or tumor)
  - Else → **VMN or VMT** (variably methylated in normal or tumor)



# Rules for DNA methylation Calls/Scores

- 5<sup>th</sup> step,
- Each gene could exist as a pair of one of the combinations
- 48 combinations could exist (3x4x4)

<b>Correlation</b>	<b>Normal Methylation</b>	<b>Tumor Methylation</b>
SNC	CUN	CUT
WNC	CMN	CMT
NNC	VMN	VMT
	IMN	IMT

# 450K (74 samples x 11,040 genes)

## Results

```
> table(genenames.450[, "N.METH"], genenames.450[, "T.METH"], genenames.450[, "correlation_call"])  
, , = NNC
```

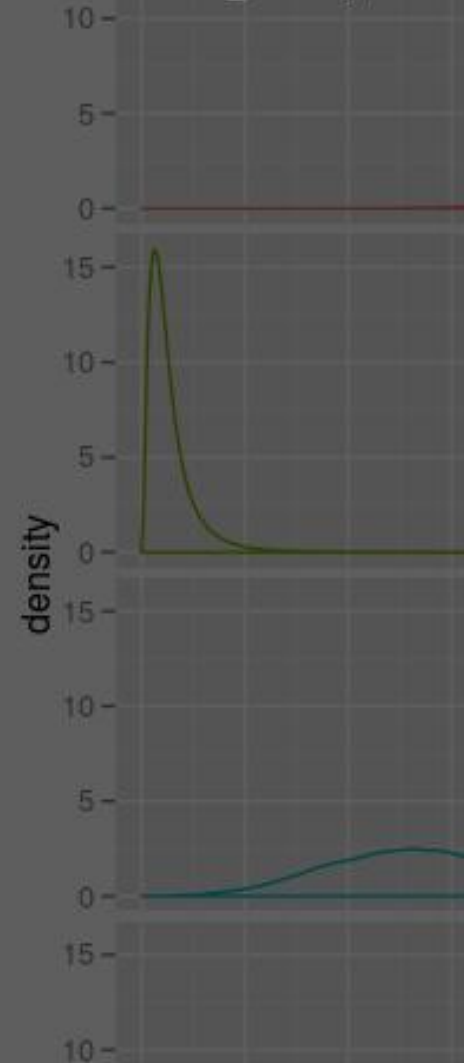
	CMT	CUT	IMT	VMT
CMN	556	0	1	240
CUN	0	781	1	97
IMN	1	0	35	97
VMN	71	14	47	597

```
, , = SNC
```

	CMT	CUT	IMT	VMT
CMN	777	0	7	268
CUN	0	1185	2	213
IMN	9	1	65	160
VMN	141	38	71	812

```
, , = WNC
```

	CMT	CUT	IMT	VMT
CMN	630	0	11	301
CUN	0	1980	3	338
IMN	2	1	60	199
VMN	111	29	95	993



# 27K (279 samples X 9,453 genes)

## Results

```
table(genenames[, "N.METH"], genenames[, "T.METH"], genenames[, "correlation_call"])  
, = NNC
```

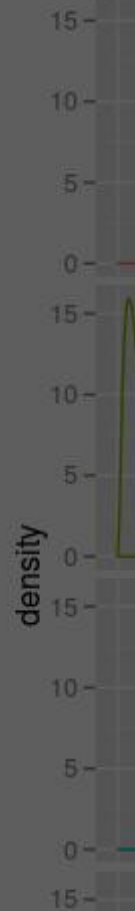
	CMT	CUT	IMT	VMT
CMN	335	0	2	432
CUN	0	4037	0	599
IMN	7	34	137	973
UNK	0	3	11	68
VMN	41	286	22	826

```
, = SNC
```

	CMT	CUT	IMT	VMT
CMN	3	0	0	7
CUN	0	15	1	44
IMN	0	1	10	66
UNK	0	0	0	0
VMN	0	3	0	29

```
, = WNC
```

	CMT	CUT	IMT	VMT
CMN	16	0	0	29
CUN	0	780	0	199
IMN	2	8	31	207
UNK	0	1	1	7
VMN	3	61	3	113



Correlation	Normal Meth	Tumor Meth
1.SNC	1.CUN	1.CUT
1.SNC	1.CUN	2.VMT
1.SNC	1.CUN	3.IMT
1.SNC	1.CUN	4.CMT
1.SNC	2.VMN	1.CUT
1.SNC	2.VMN	2.VMT
1.SNC	2.VMN	3.IMT
1.SNC	2.VMN	4.CMT
1.SNC	3.IMN	1.CUT
1.SNC	3.IMN	2.VMT
1.SNC	3.IMN	3.IMT
1.SNC	3.IMN	4.CMT
1.SNC	4.CMN	1.CUT
1.SNC	4.CMN	2.VMT
1.SNC	4.CMN	3.IMT
1.SNC	4.CMN	4.CMT
2.WNC	1.CUN	1.CUT
2.WNC	1.CUN	2.VMT
2.WNC	1.CUN	3.IMT
2.WNC	1.CUN	4.CMT
2.WNC	2.VMN	1.CUT
2.WNC	2.VMN	2.VMT
2.WNC	2.VMN	3.IMT
2.WNC	2.VMN	4.CMT
2.WNC	3.IMN	1.CUT
2.WNC	3.IMN	2.VMT
2.WNC	3.IMN	3.IMT
2.WNC	3.IMN	4.CMT
2.WNC	4.CMN	1.CUT
2.WNC	4.CMN	2.VMT
2.WNC	4.CMN	3.IMT
2.WNC	4.CMN	4.CMT
3.NNC	1.CUN	1.CUT
3.NNC	1.CUN	2.VMT
3.NNC	1.CUN	3.IMT
3.NNC	1.CUN	4.CMT
3.NNC	2.VMN	1.CUT
3.NNC	2.VMN	2.VMT
3.NNC	2.VMN	3.IMT
3.NNC	2.VMN	4.CMT
3.NNC	3.IMN	1.CUT
3.NNC	3.IMN	2.VMT
3.NNC	3.IMN	3.IMT
3.NNC	3.IMN	4.CMT
3.NNC	4.CMN	1.CUT
3.NNC	4.CMN	2.VMT
3.NNC	4.CMN	3.IMT
3.NNC	4.CMN	4.CMT

# Rules for DNA methylation Calls/Scores

- 5<sup>th</sup> step cont.
- Need to decide appropriate “call” and “score” labels.

Call	Desc
<b>MG</b>	Methylation gain compared to normal
<b>ML</b>	Methylation loss compared to normal
<b>MT</b>	Methylated in tumor
<b>UT</b>	Unmethylated in tumor
<b>ES</b>	Epigenetically silenced
<b>UC</b>	Unable to make call

Score	Desc (confidence)
<b>4</b>	High
<b>3</b>	Med-High
<b>2</b>	Med
<b>1</b>	Low
<b>0</b>	No call



Correlation	Normal Meth	Tumor Meth
1.SNC	1.CUN	1.CUT
1.SNC	1.CUN	2.VMT
1.SNC	1.CUN	3.IMT
1.SNC	1.CUN	4.CMT
1.SNC	2.VMN	1.CUT
1.SNC	2.VMN	2.VMT
1.SNC	2.VMN	3.IMT
1.SNC	2.VMN	4.CMT
1.SNC	3.IMN	1.CUT
1.SNC	3.IMN	2.VMT
1.SNC	3.IMN	3.IMT
1.SNC	3.IMN	4.CMT
1.SNC	4.CMN	1.CUT
1.SNC	4.CMN	2.VMT
1.SNC	4.CMN	3.IMT
1.SNC	4.CMN	4.CMT
2.WNC	1.CUN	1.CUT
2.WNC	1.CUN	2.VMT
2.WNC	1.CUN	3.IMT
2.WNC	1.CUN	4.CMT
2.WNC	2.VMN	1.CUT
2.WNC	2.VMN	2.VMT
2.WNC	2.VMN	3.IMT
2.WNC	2.VMN	4.CMT
2.WNC	3.IMN	1.CUT
2.WNC	3.IMN	2.VMT
2.WNC	3.IMN	3.IMT
2.WNC	3.IMN	4.CMT
2.WNC	4.CMN	1.CUT
2.WNC	4.CMN	2.VMT
2.WNC	4.CMN	3.IMT
2.WNC	4.CMN	4.CMT
3.NNC	1.CUN	1.CUT
3.NNC	1.CUN	2.VMT
3.NNC	1.CUN	3.IMT
3.NNC	1.CUN	4.CMT
3.NNC	2.VMN	1.CUT
3.NNC	2.VMN	2.VMT
3.NNC	2.VMN	3.IMT
3.NNC	2.VMN	4.CMT
3.NNC	3.IMN	1.CUT
3.NNC	3.IMN	2.VMT
3.NNC	3.IMN	3.IMT
3.NNC	3.IMN	4.CMT
3.NNC	4.CMN	1.CUT
3.NNC	4.CMN	2.VMT
3.NNC	4.CMN	3.IMT
3.NNC	4.CMN	4.CMT

Call if TM < 0.25	Confidence Score if TM < 0.25
UT	4
UT	3
UT	3
UT	3
ML	1
UT	3
UT	3
UT	2
ML	2
UT	3
UT	3
UT	1
ML	4
ML	4
ML	3
ML	1
UT	4
UT	3
UT	3
UT	3
ML	1
UT	3
UT	3
UT	3
ML	2
UT	3
UT	3
UT	3
UT	1
ML	4
ML	4
ML	3
ML	1
UT	4
UT	3
UT	3
UT	3
ML	1
UT	3
UT	3
UT	3
ML	2
UT	3
UT	3
ML	4
ML	4
ML	3
ML	1

Call if 0.25 < TM < 0.75	Confidence Score if 0.25 < TM < 0.75
ES	1
ES	3
ES	2
ES	3
UC	0
UC	0
UC	0
UC	0
UC	0
UC	0
UC	0
ML	3
ML	3
ML	2
ML	2
MG	2
ES	2
ES	1
ES	2
UC	0
UC	0
UC	0
UC	0
UC	0
UC	0
UC	0
ML	3
ML	3
ML	2
ML	2
MG	2
MG	3
MG	2
MG	3
UC	0
UC	0
UC	0
UC	0
UC	0
UC	0
UC	0
ML	3
ML	3
ML	2
ML	2

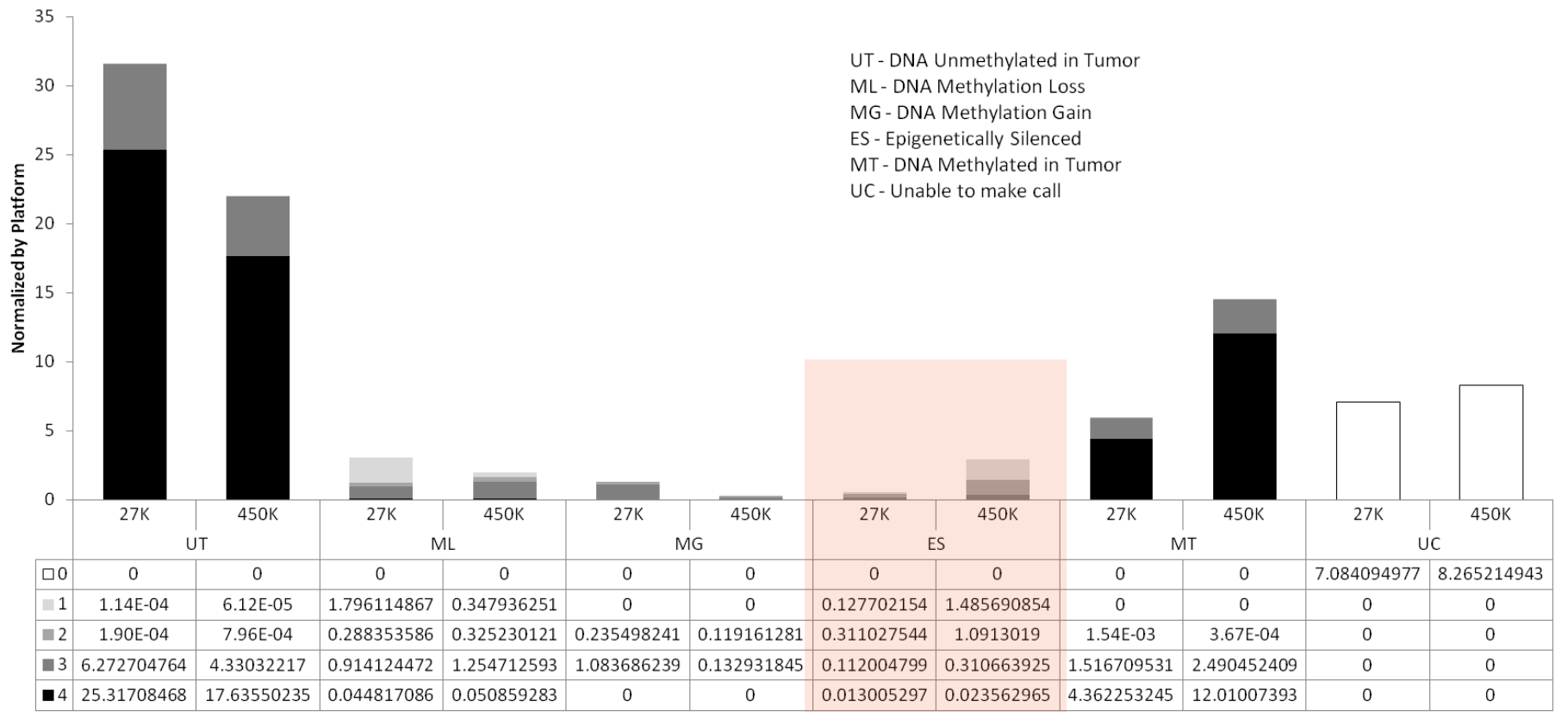
Call if TM > 0.75	Confidence Score if TM > 0.75
ES	2
ES	4
ES	3
ES	4
ES	1
ES	1
ES	2
ES	1
ES	1
ES	1
ES	2
MT	2
MT	3
MT	3
MT	4
MG	2
ES	3
ES	2
ES	3
MT	2
MT	4
MT	4
MT	4
MT	2
MT	4
MT	4
MT	4
MT	2
MT	3
MT	3
MT	4
MG	2
MG	3
MG	3
MG	4
MT	2
MT	4
MT	4
MT	4
MT	2
MT	4
MT	4
MT	4
MT	2
MT	3
MT	3
MT	4

# Rules for DNA methylation Calls/Scores

- 6<sup>th</sup> step,
- Using the table rule (from step 5) as reference, we made “calls” on each gene per sample, see R script for exact steps.
- Basically, each cell in the DNA methylation data matrix (samps X genes) is evaluated.
  - Beta-values  $< 0.25$
  - Beta-values  $> 0.25$  &  $< 0.75$
  - Beta-values  $> 0.75$

27K: 279 samp X 9453 (x2) = 5,274,774  
 450K: 74 samp X 11040 (x2) = 1,633,920

**Distribution of DNA methylation calls and scores  
 across 27K (n=279) & 450K (n=74)**



# Data and code available in dropbox

- Finally,
- Merged 27K/450K based on geneID.
- Added data/codes etc. to dropbox.
- See following path:  
“Dropbox/TCGA\_GBM\_MS\_Writing\_Group/Working\_Group\_SampleManifest/DNA\_methylation”
- R file (which contains the entire data matrix with calls, scores), R scripts (most of the code is listed), exported data as tab-delimited file (for those non R users). “Calls” and “Scores” are separated in two data matrix.
- R file:  
“TCGA\_GBM\_DNAMETHYLATION\_CALLS\_SCORES\_20120112\_Noushmehr\_ver2.rda”