# What's new
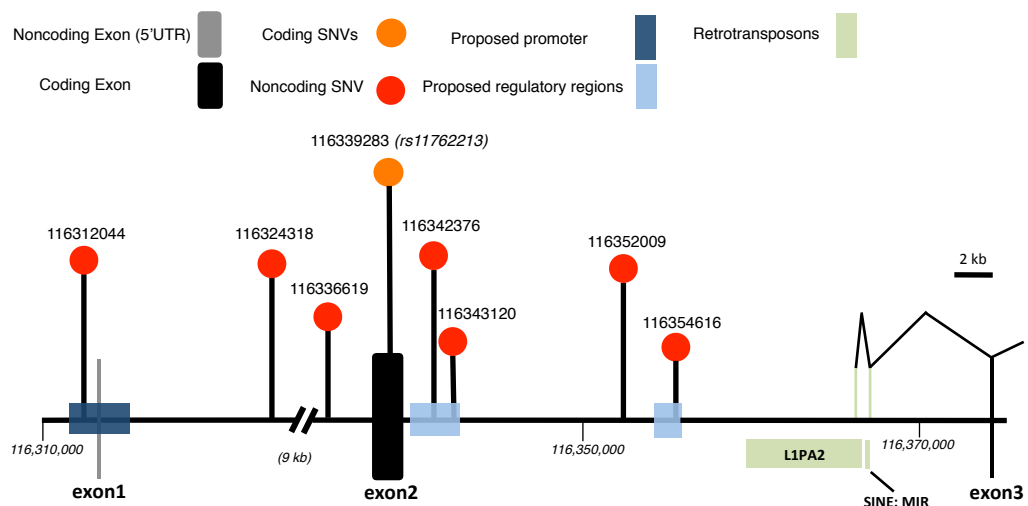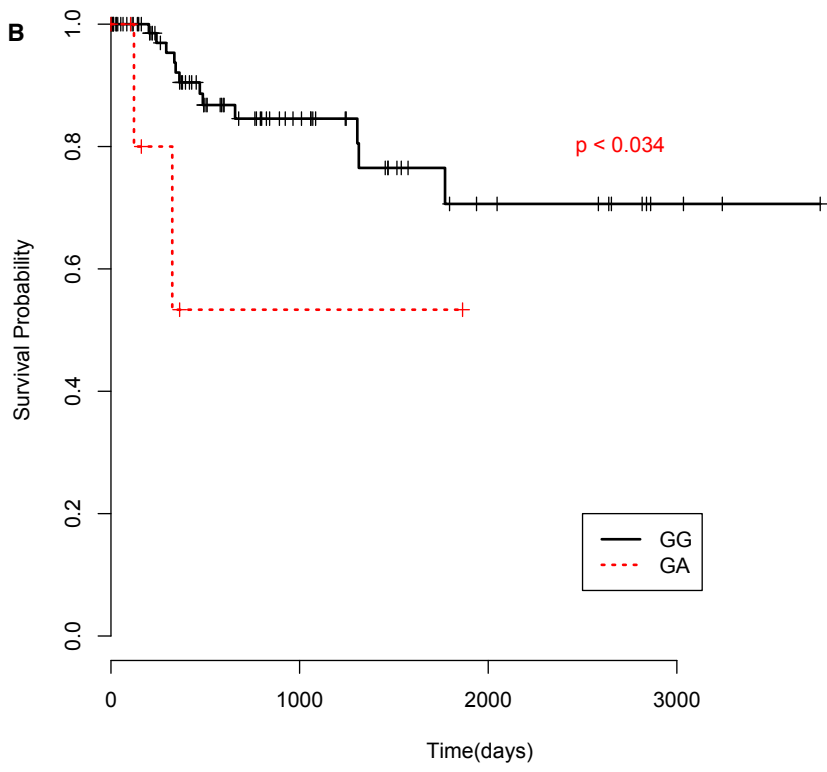
- Figure re-plot
- SV comparison
- New DHS/RepliSEQ
- Evolution tree
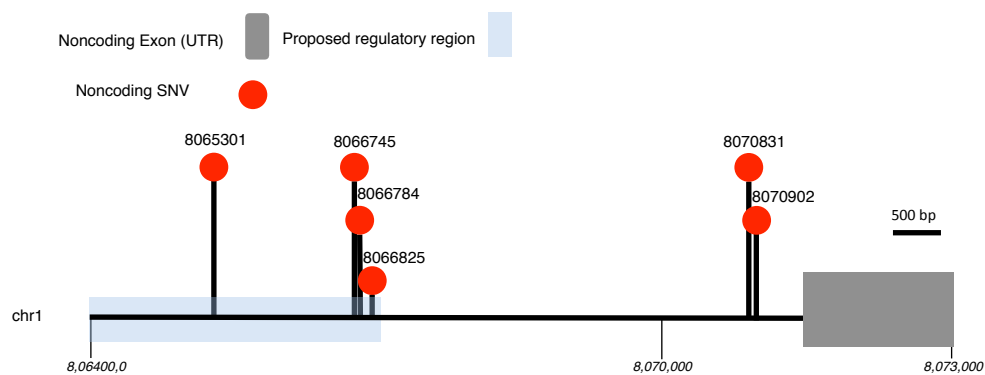
**A** MET

Noncoding Exon (5'UTR)
Coding Exon
Coding SNVs
Noncoding SNV
Proposed promoter
Proposed regulatory regions
Retrotransposons

116339283 *(rs11762213)*

116312044
116324318
116336619
116342376
116343120
116352009
116354616

chr7

*116,310,000*
*(9 kb)*
*116,350,000*
*116,370,000*

2 kb

L1PA2
SINE: MIR

**exon1**
**exon2**
**exon3**

**B**

Survival Probability

p < 0.034

GG
GA

Time(days)

**A** ERRFI1

Noncoding Exon (UTR) ▮ Proposed regulatory region ▮

Noncoding SNV ⬤

8065301  8066745  8070831

8066784  8070902

8066825

500 bp

chr1

8,064000,0  8,070,000  8,073,000

**B** NEAT1

Noncoding RNA (NEAT1) ▮ Proposed promoter ▮

Noncoding SNV ⬤ Proposed regulatory region ▮

65192209  65195872  65199075

65194037  65196991  65200157

65201466

1 kb

chr11

65,190, 000  65,200,000  65,210,000  65,214,000

**C** p<0.032

Normalized counts

60000

20000

0

wt.  mut.

**D**

Survival Probability

0.9

0.7

0.5

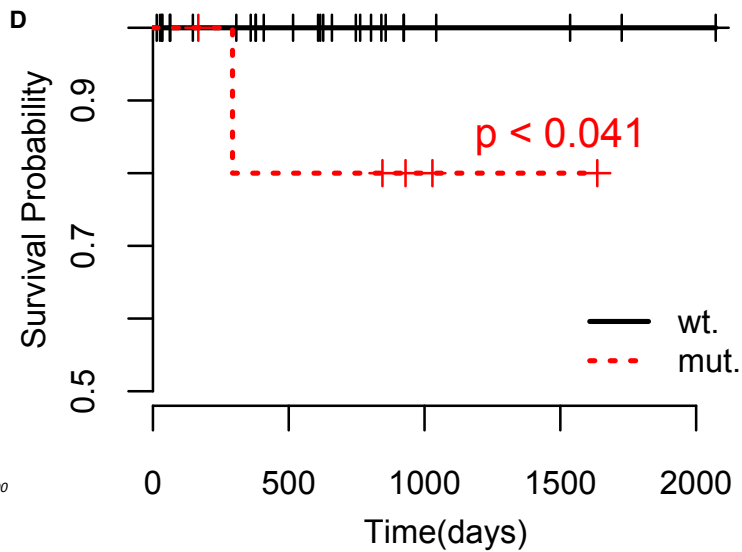p < 0.041

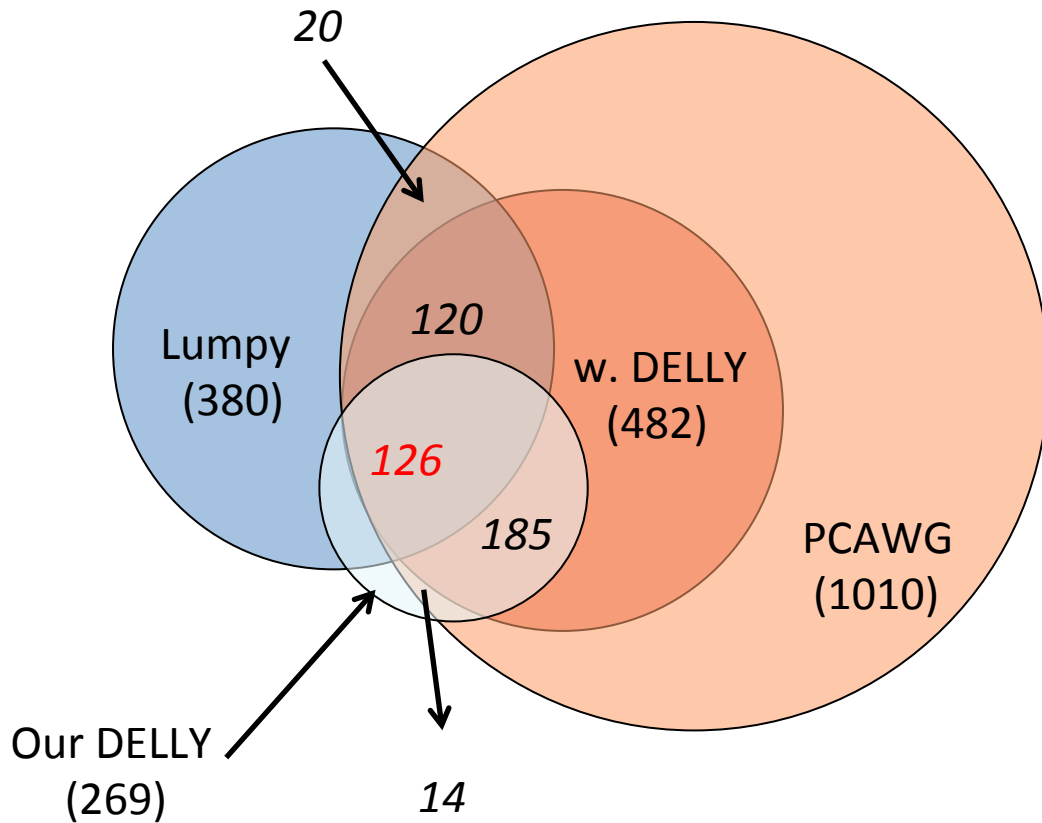wt.
mut.

0  500  1000  1500  2000

Time(days)

# SV comparison

- First, why Lumpy/SVscore?
  - Our strength at bkpts (as to aCGH)
    - Overlapping bkpts with functional regions
    - What about bkpt labeled as IMPRECISE
  - Lumpy/SVscore gives the right normalization
    - Normalized CADD scores around bkpts
    - Report intervals as well
  - Issues: Quality & SVTyper
    - The genotyper works weirdly

# SV comparison



This is on 32 samples
(*we called SV from 35*)

Criteria:
1. *0.5* reciprocal overlap
2. Matched sample
3. Matched SV CLASS

Only get *126* SVs if overlapping DELLY with Lumpy using 32 samples

# SV comparison

- If we think PCAWG DELLY is the ground truth
  - Other methods rely on assembly…different searching space (complicated SV events produce way too many bkpts)
  - Not a fair comparison

| | #SVs | FN | FP | Addl. catch |
|---|---|---|---|---|
| Lumpy | 380 | 71% | 63.2% | 20 (5.26%) |
| Delly (ours) | 269 | 61.6% | 31.2% | 14 (5.20%) |

# Bkpts&interval overlapping

- No overlapping with pRCC MutSig genes (~10)
- COSMIC genes

- Three overlapping with pRCC MutSig
  - Extensive SVs in TCGA-B9-4116
  - DEL of STAG2
  - A large INV involves NFE2L2

|   | bkpts |
|---|-------|
| 1 | AFF3 |
| 1 | AKAP9 |
| 1 | ATRX |
| 1 | CDK6 |
| 1 | CDKN2A |
| 1 | CLP1 |
| 1 | CREBBP |
| 1 | ERCC2 |
| 1 | LRP1B |
| 1 | MKL1 |
| 1 | MLLT10 |
| 1 | POLE |
| 1 | POT1 |
| 1 | SND1 |
| 1 | SOX2-OT |
| 1 | STAT5B |
| 1 | STK11 |
| 1 | THRAP3 |
| 2 | SPEN |

# SV overlapping

- Some interesting COSMIC cases
  - 3 CDKN3A (confirmed NEJM 3/5 cases)

  - 2 SDHB deletion (interacts with FH and SDHA)
  - 1 EGFR duplication (pRCC responses to TKIs)
  - 1 HIF1A duplication (inhibited by VHL;)
  - 1 polyE bkpts (same case)
  - 1 DNMT3A deletion (affect methylation)

  - 1 MALAT1 *deletion*
  - 1 HGF *deletion (ligand for MET)*

| | |
|---|---|
| 1 | POLE2 |
| 1 | DNMT3A |
| 1 | HIF1A |
| 1 | IDH1 |
| 1 | MALAT1 |
| 2 | EGFR |
| 2 | SDHB |
| 3 | CDKN2A |
| | |
| 1 | HGF |

# New DHS scheme

- Pulled fetal kidney cortex DHS from Roadmap

- 11 samples from different fetuses
  - Ultra-conservative: take the overlap of all samples
  - DHS percentage
    roughly matches ENCODE
    HEK293 narrowpeak
    (previously used)

# RepliSEQ

- Newly defined CR-genes
  - Overlapping NEJM spreadsheet(pRCC-related) genes with CR&SWI/SNF pathway gene list
- Conservative RT signals (*Ref: NGen.*)
  - Taken median from 11 ENCODE cell lines
- Prudent adapted KS test (subsampling)
  - Q: does the RT dist. of SNVs from CR$^{mut.}$ samples differ from CR$^{wt.}$ samples
  - Randomly shuffle labels (9 v.s. 26) for 1,000 times to generate imperial TS distribution

# Why we need subsampling



~60% sampling produce
p < 0.0001

# KS tests w/ subsampling: Issues

- A. p-value below limits
  - Our test produces p < 2.2e-16 (which becomes 0)
  - However, we observed 26% such cases
- B. test power not uniform
  - SNVs in each sample varies
  - Makes test statistic (D) not directly comparable
  - Worse: power is not uniform!

# Evolution tree



- Less interesting because:
  - Single punch from each patient
  - Not ultra-deep sequence
  - Masked all CNV + MAF>0.6

TCGA-A4-A48D

KDM5A:nonsynonymous
KDM5A:nonsynonymous

0.52
0.84

TCGA-A4-A4ZT

0.78 0.4

TCGA-A4-A57E

0.83 0.44

TCGA-AL-3466

NEAT1:nc:hotspot

0.75 0.37

TCGA-AL-3468

TRIM33:nonsynonymous
NEAT1:nc:hotspot
DCTN1:nonsynonymous
CRTC2:nonsynonymous
TRIP11:nonsynonymous

0.88 0.24

TCGA-AL-3472

ERRFI1:nc:hotspot
SF3B1:nonsynonymous

0.87 0.35

TCGA-AL-3473

BIRC3:nonsynonymous
NSD1:nonsynonymous

0.35

TCGA-AL-A5DJ

KMT2C:nonsynonymous
ARID1A:nonsynonymous

0.52
0.87
0.34

TCGA-B1-A47M

NFE3L2:nonsynonymous
NF2:nonsynonymous

0.72 0.37

TCGA-B1-A47N

0.83

TCGA-B1-A47O

HERPUD1:nonsynonymous

0.48

TCGA-B3-3925

ROS1:nonsynonymous
NEAT1:nc:hotspot
NEAT1:nc:hotspot

0.72 0.31

TCGA-B3-3926

ELF4:nonsynonymous

0.98 0.59 0.39

TCGA-B3-4103

NEAT1:nc:hotspot
FANCG:nonsynonymous

0.78 0.73

TCGA-B9-4113

ERRFI1:nc:hotspot
LIFR:nonsynonymous

0.75 0.54
0.36

TCGA-B9-4114

KMT2C:nonsynonymous
KTP1A:nonsynonymous

0.61
0.93

FALB2:nonsynonymous
ROS1:nonsynonymous

TCGA-B9-4115

0.75 0.68

TCGA-B9-4116

PBX1:nonsynonymous
MYCN:nonsynonymous
EP300:nonsynonymous
NUP214:nonsynonymous
FAT1:nonsynonymous

0.76 0.42

TCGA-B9-4117

HERPUD1:nonsynonymous
CREBBP:nonsynonymous

0.98 0.62 0.35

TCGA-B9-4617

0.51
0.89

TCGA-B9-A44B

FTP4B:nonsynonymous
SRSN:nonsynonymous
EIF4A2:nonsynonymous
KA7BB:nonsynonymous

1 0.85 0.47

TCGA-GL-A4EM

ATR:nonsynonymous
MET:nonsynonymous

TCGA-GL-A59R

XPO1:nonsynonymous
JAK2:nonsynonymous
SF3B1:nonsynonymous
GPHN:nonsynonymous

0.84 0.35
0.49

TCGA-HE-A5NF

0.56
0.84

TCGA-HE-A5NH

NUP214:nonsynonymous

0.81
0.46

TCGA-HE-A5NI

NIN:nonsynonymous
EPS15:nonsynonymous
ERRFI1:nc:hotspot

0.77 0.41

TCGA-HE-A5NJ

0.83

TCGA-HE-A5NL

NCOA1:nonsynonymous
NEAT1:nc:hotspot

0.84

TCGA-IA-A40X

IL21R:nonsynonymous
IL21R:nonsynonymous
NUP98:nonsynonymous

0.79 0.34

TCGA-IA-A40Y

0.53
0.89

FGFR2:nonsynonymous
SF3B1:nonsynonymous

0.37

TCGA-MH-A55W

ELK4:removedStop

0.99 0.71

TCGA-MH-A55Z

SMARCA4:nonsynonymous
NEAT1:nc:hotspot
DNMT3A:prematureStop
RANBP2:nonsynonymous

0.85

TCGA-MH-A560

0.74

TCGA-MH-A561

MET:nonsynonymous
GPHN:nonsynonymous
NRG1:nonsynonymous
SS18L1:nonsynonymous
NCOA2:nonsynonymous
NEAT1:nc:hotspot

0.79

TCGA-MH-A562

ERRFI1:nc:hotspot
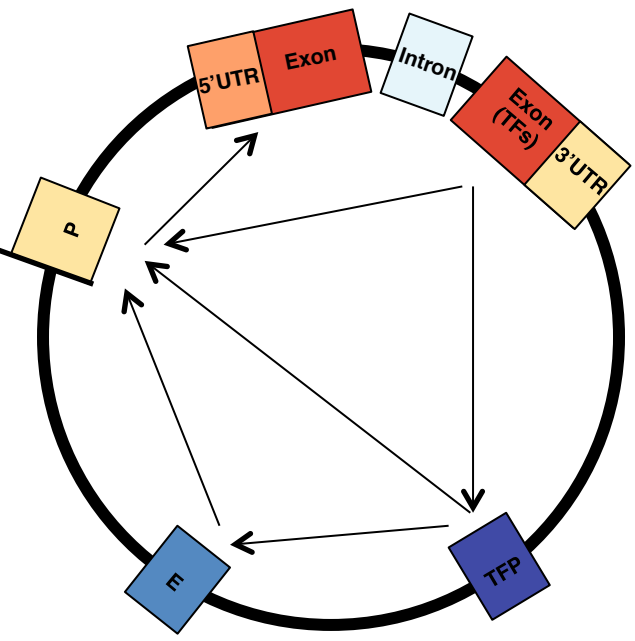FAT1:nonsynonymous
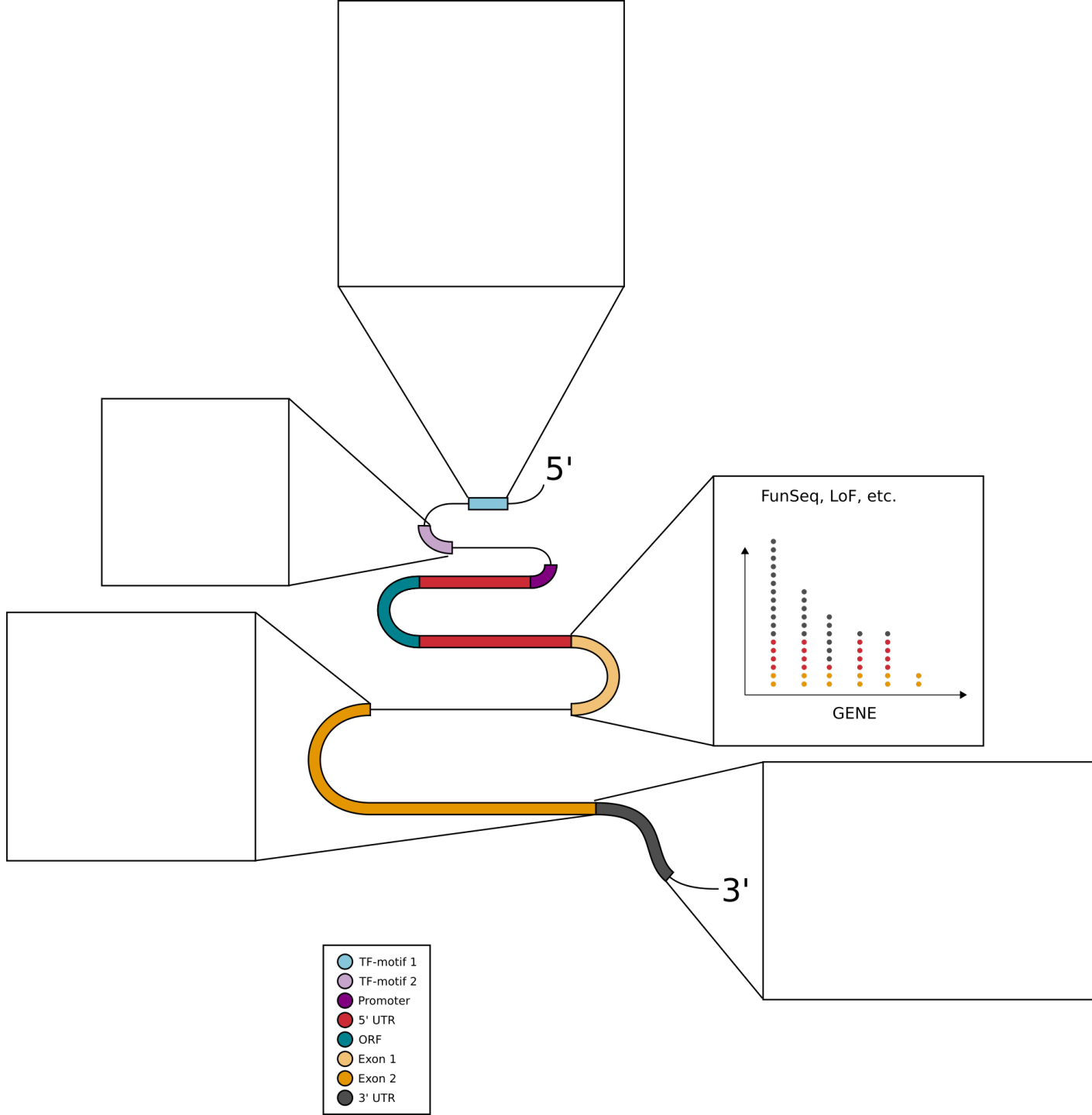SLC34A2:nonsynonymous

0.71 0.35

# Paper E thoughts

- What we want to sell
  - A complete "tumor portrait"
  - Completeness on several levels
    - Unprecedentedly large WGS cohort
    - Coding/noncoding/SVs…

- Fig1: a fancy figure to visualize the deluge of the data
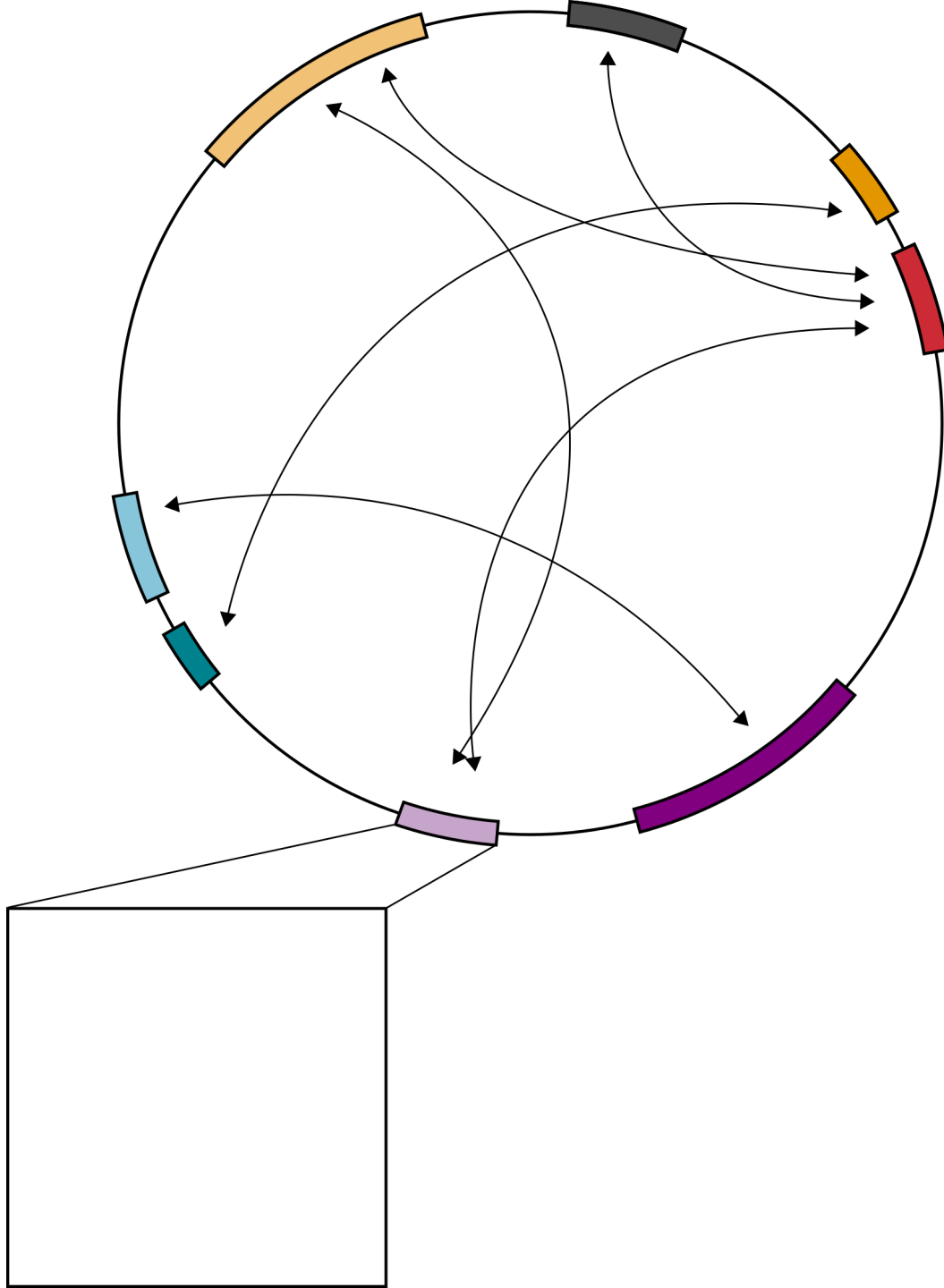  - Show a gene ("genecreek" or "-circle")

# Issues

- One sample with 23 EPB41L2 promoter mutation?
  - And 9 has the identical FunSeq score!
  - Annotation issue? This promoter spans >170kb!

  - However, this is an excluded sample...

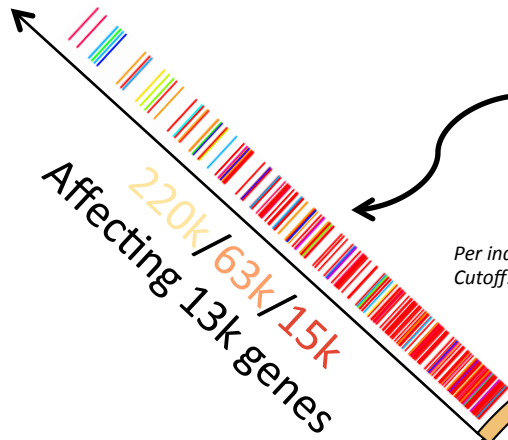  - Still, another included sample has 8...

325k/92k/22k
Affecting 14.2k genes

5'UTR

Exon

Intron

Exon (TFS)

3'UTR

P

TFP

E

5'

3'

FunSeq, LoF, etc.

GENE

TF-motif 1
TF-motif 2
Promoter
5' UTR
ORF
Exon 1
Exon 2
3' UTR

Gene functions? FunSeq breakdowns?

Affecting 13k genes

220k / 63k / 15k

Per indiv. per gene (57k)
Cutoff: 1.5

| 108 | TCF4 |
| 109 | DUSP22 |
| 124 | TERT |
| 125 | SEPT9 |
| 154 | ZNF595 |

# Annotation

```
   51 volgesteinDriver.oncogene
   58 volgesteinDriver.TSG
  101 cellCycleGene.
  114 apoptosisGene.
  122 dnaRepairGene.
  339 cancerPathwayGene.
  385
 1375 Metabolic_Genes_RCT
 2646 essentialGene.
 2905 immuneResponseGene.
 9969 nonEssentialGene.
```

1. Normalize by the #gene
2. Normalize by the length of elements
3. Correct for RT/Trinucleotide context

# Some ideas