

# Recurrence analysis considering sequence composition and variable mutation rates

Inigo Martincorena  
ICGC pancancer phone call  
17<sup>th</sup> Nov 2014

# Factors to consider

## 1. Length of the element



# Factors to consider


## 2. Trinucleotide mutation rates and sequence composition

E.g.: C>T at CpGs tend to have a rate per bp  $\sim 40x$  higher than other substitutions

*Random mutations:*

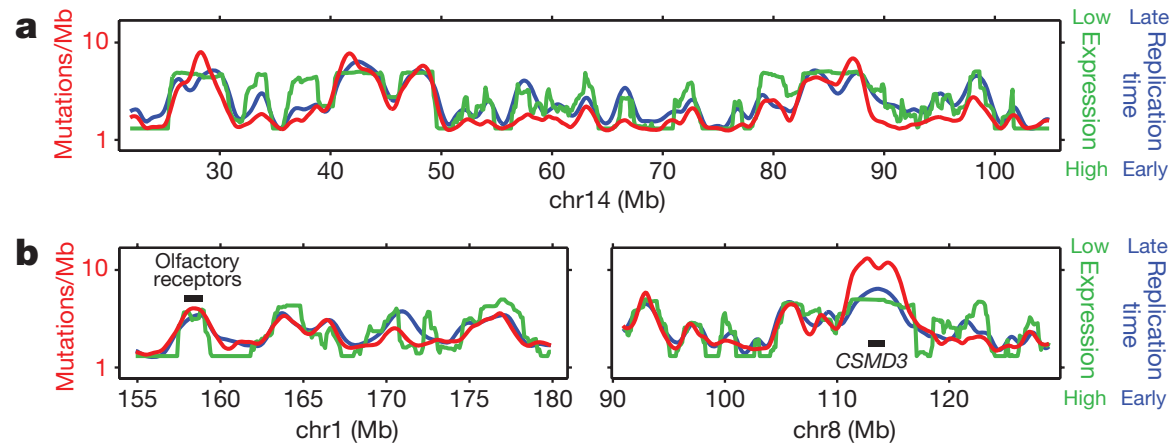


 *Element 1*

 *Element 2 (e.g. TFBS, phosphorylation site...)  
could be significant with uniform models*

# Factors to consider

## 3. Regional variation of the mutation rate



### Possible solutions:

- Use local density of mutations (e.g. silent sites, 10kb around...)
- Use covariates (e.g. expression, repl time...)

*But uncertainty in the estimates needs to be included in the calculation of p-values.*

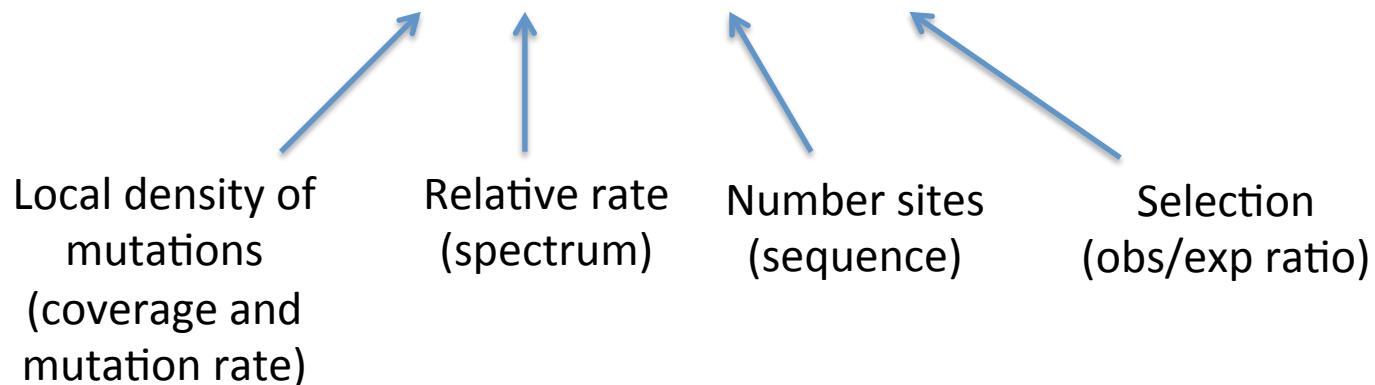


# Our approach on genes

*e.g.* frequency of synonymous and missense A>C mutations

$$\lambda_{syn,AtoC} = (t) * (AtoC) * (L_{syn,AtoC})$$

$$\lambda_{mis,AtoC} = (t) * (AtoC) * (L_{mis,AtoC}) * (wMIS)$$



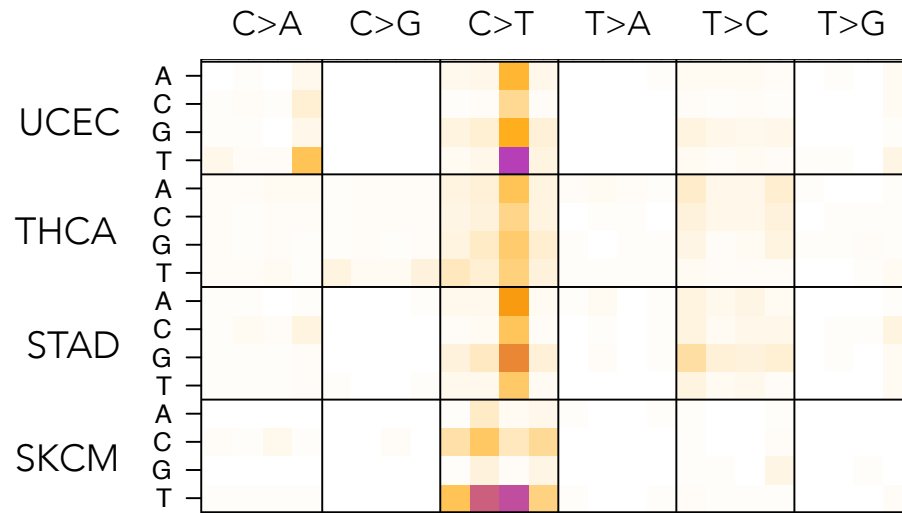
# Adding context-dependence

Context effects easily incorporated as multiplicative rates

$$\lambda_{mis,CtoT} = (t)*(CtoT)*(L_{mis,CtoT})*(wMIS) \quad 12 \text{ rates model}$$

$$\lambda_{mis,CtoT} = (t)*(CtoT)*(L_{mis,CtoT})*(wMIS)*(CpGctx)*(TpCctx) \quad 16 \text{ rates model}$$

$$\lambda_{mis,CtoT} = (t)*(ACGtoATG)*(L_{mis,ACGtoATG}) \quad 192 \text{ rates model}$$



# Local mutation rates



*False positives*

*Low sensitivity*

**Model 1:**  $\mathbf{bg = t}$

Assumes no variation between genes.  
Gives hundreds of FP in large datasets.

**Model 2:**  $\mathbf{bg = Pois(t | n_{syn})}$

Very conservative.  
Clean results for large datasets.  
It does not exploit info from other genes.



# Local and global model

## *Exploiting covariates: negative binomial regression ( $t * trinuc\_rates * L$ )*

- Trinucleotide rates and sequence composition included as the offset (fixed factors)
- Negative binomial: mutations are Poisson, rates among genes Gamma distributed
- Adding covariates reduces the uncertainty (dispersion) of the background  $\rightarrow$  higher power

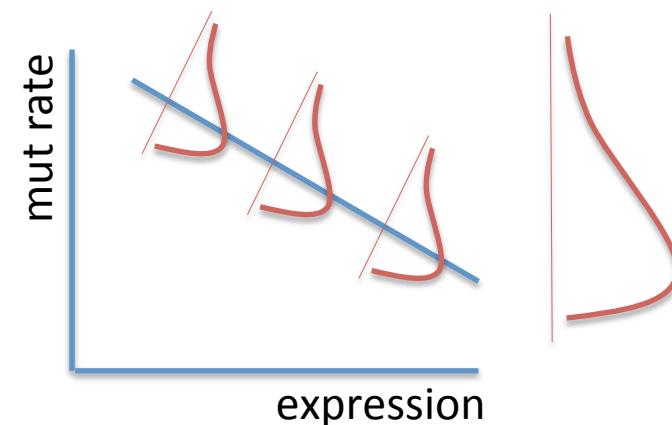
*Covariates are fitted to each dataset, not assumed*

### **Negative binomial regression**

`glm.nb(obs ~ offset(expected) + expr + reptime ...)`

	Obs_syn	Exp_syn	Expression	Repl time
Gene A	10	8.53	20012	4.35
Gene B	38	30.76	1782	2.67
Gene C	1	2.45	59927	-8.6
...	...	...	...	...
Gene Z	290	263	306	4.56

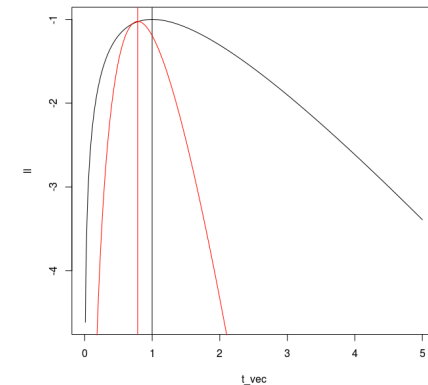
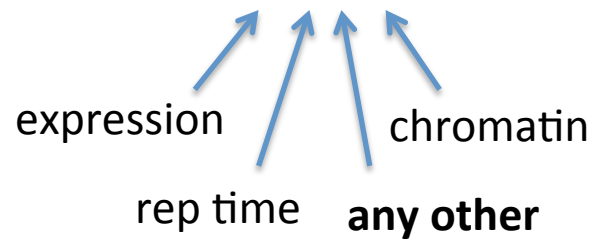
**bg = Gamma( $t$  | covariates)**



# Local mutation rates



$$bg = \text{Pois}(t | n_{\text{syn}}) * \text{Gamma}(t | \text{covariates})$$



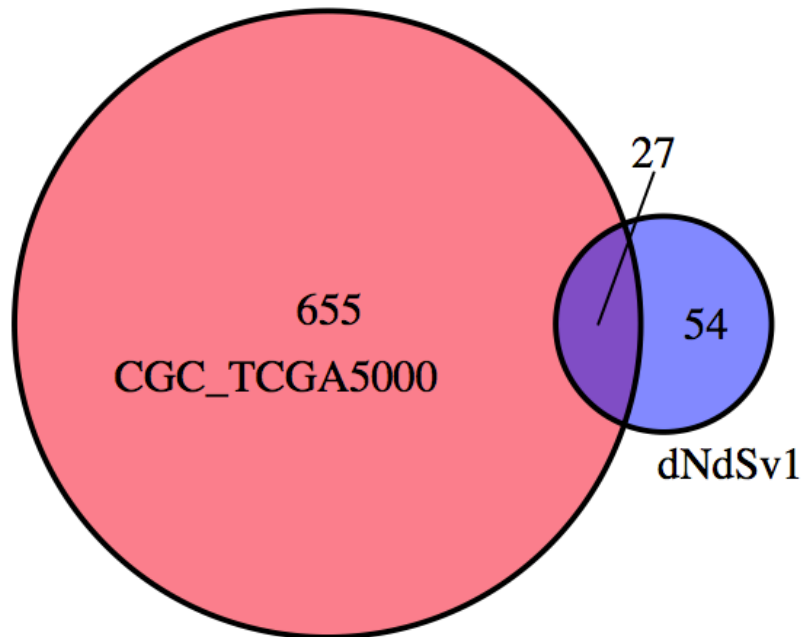
*Model 3*

$$bg = \text{Pois}(t | n_{\text{syn}}) * \text{Gamma}(t | \text{cov})$$

$$bg = \text{Pois}(t | n_{\text{syn}}) * \text{Gamma}(t | \text{cov})$$

# 607 pilot (gene analysis only)

## A) Only correcting for sequence composition



**Overlap = 33.3%**

Dubious novel cancer genes as the list includes TTN and other large genes as well as interesting candidates.

**Random dataset:**

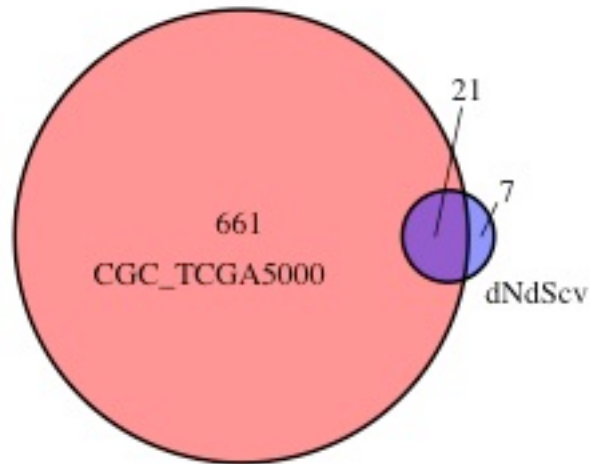
9 false positives

**Lowest qval = 4e-9**

Obviously the background model is not realistic

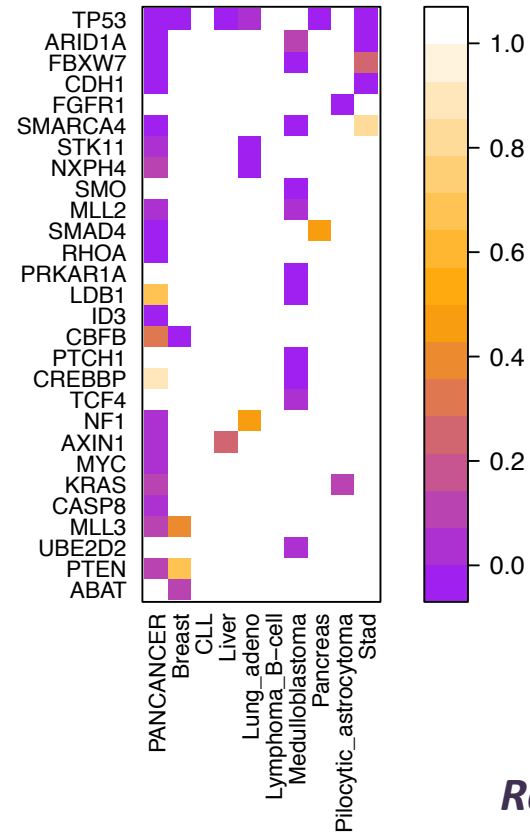
# 607 pilot (gene analysis only)

## B) Accounting for sequence composition and variable mutation rates across genes



**Overlap = 75%**

Good novel cancer genes: At least 5 of the 7 new candidates are known to be involved in cancer (PTCH1, ID3, LDB1, TCF4 and ABAT).



**Random dataset shows very good behaviour**

0 false positives  
Lowest qual = 0.86

# 607 pilot (gene analysis only)

## Some conclusions and ideas for noncoding sites:

- Need to account for extreme trinucleotide rates: especially for hotspot analysis and recurrence in small motifs (e.g. TFBS, phosphosites, functional bias...)
- Need to account for variable mutation rates
  - *Local estimates and/or covariates*
  - *And explicitly consider the uncertainty that remains in the model*
- Increasing power within this framework:
  - *Site recurrence or local hotspots (same model applies to a single site)*
  - *Classify sites by importance: nonsense vs missense, polyphen, TFBS within promoters...*
- FDR correction has to be decided *a-priori* (e.g. correct for all genes/sites even if p-values were calculated only for those than had mutations)