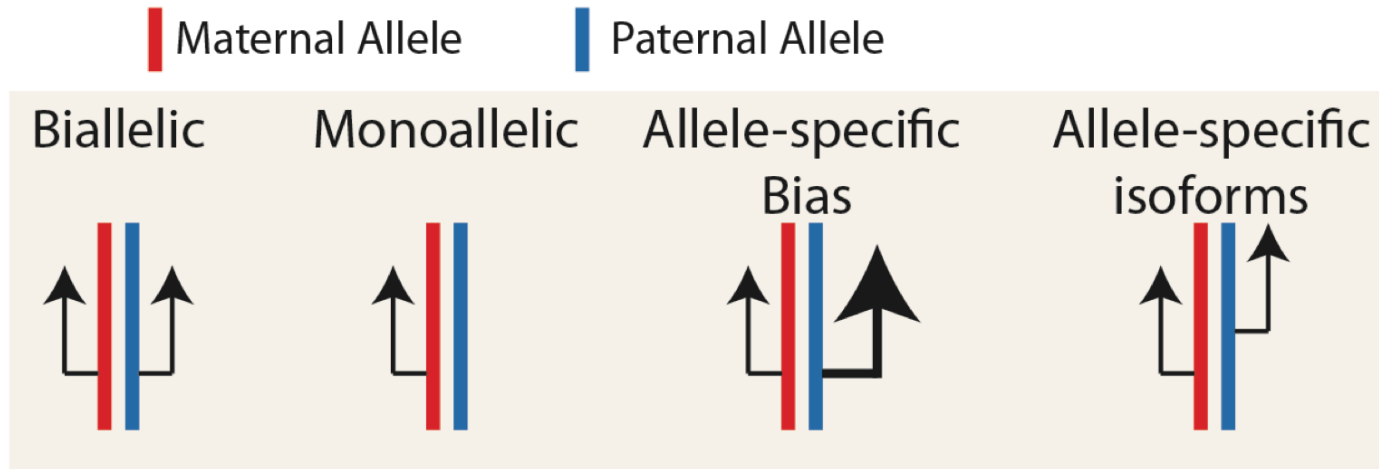


Studying allele specific expression to capture condition-specific genetic regulatory effects

# Overview

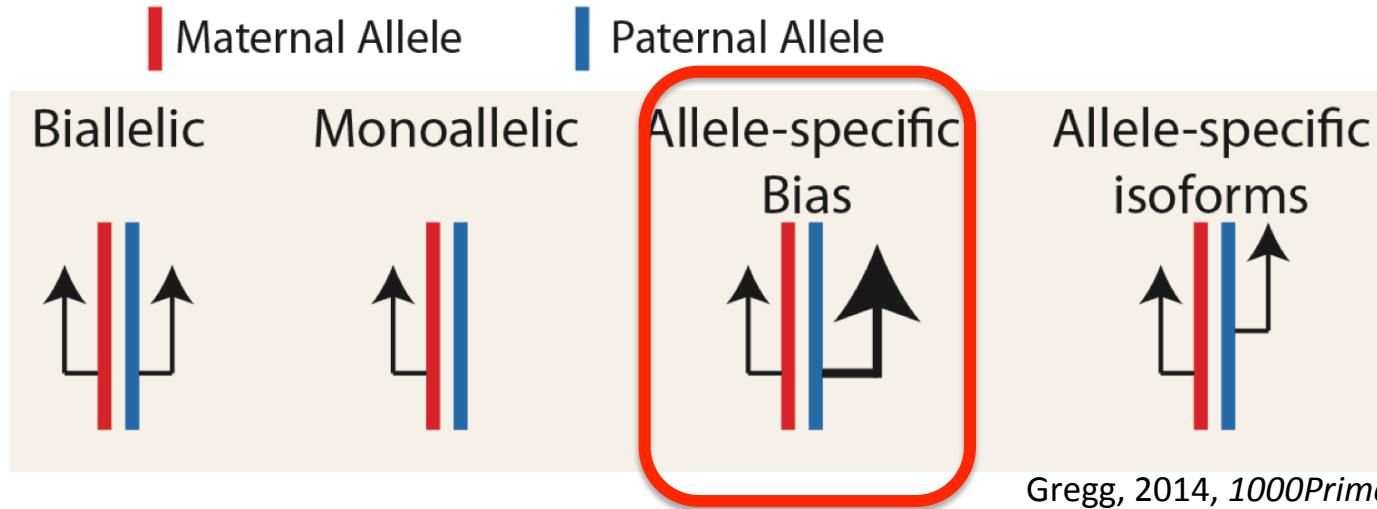
- Quick intro to allele specific expression (again)
- New stuff in pipeline
- TASC ASE
- HLA class II
- STAT pilot

# Types of allelic expression

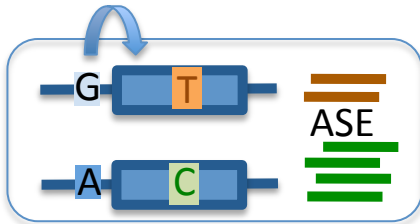


Gregg, 2014, *1000Prime reports*

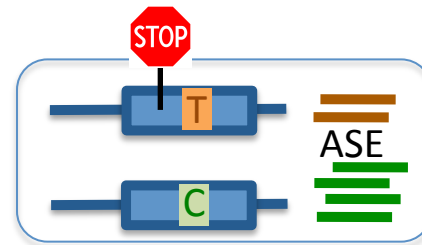
# Types of allelic expression



Genetic regulatory variation

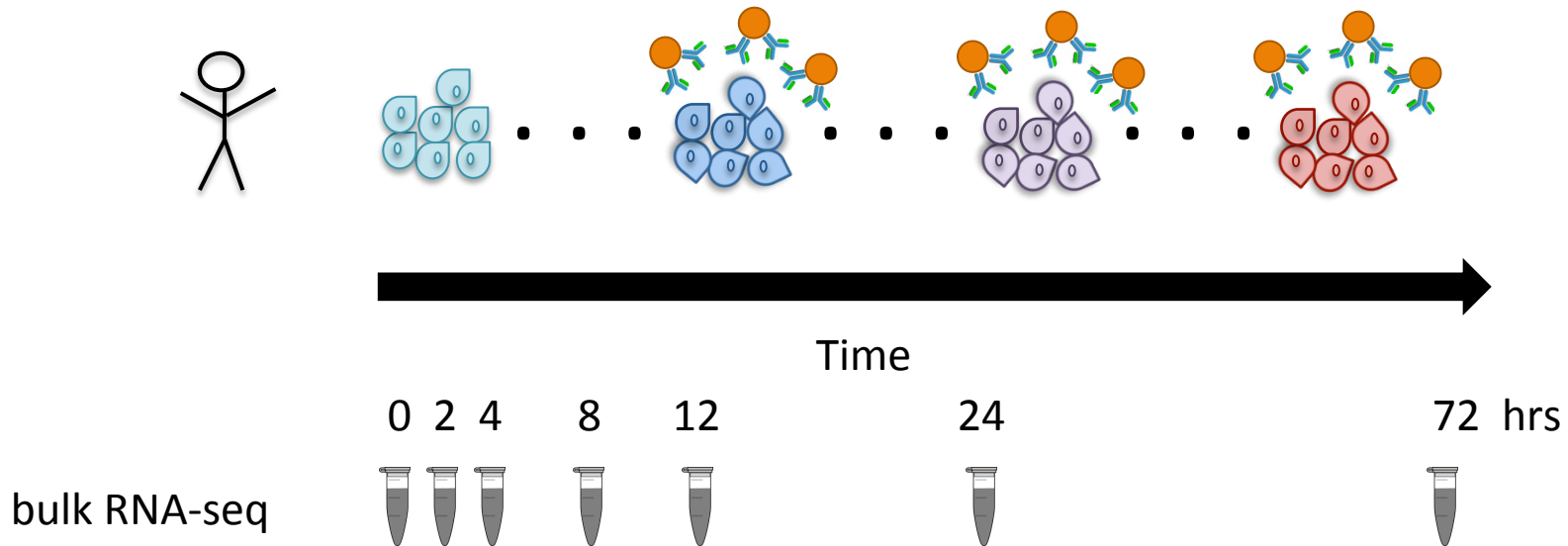


Nonsense mediated decay



# ASE in response to stimuli

CD4<sup>+</sup> T<sub>MEM</sub> cells stimulated with anti-CD3/CD28 beads

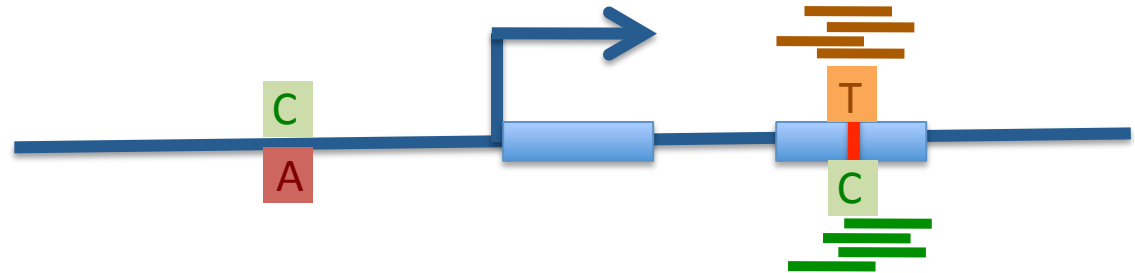


## Objective

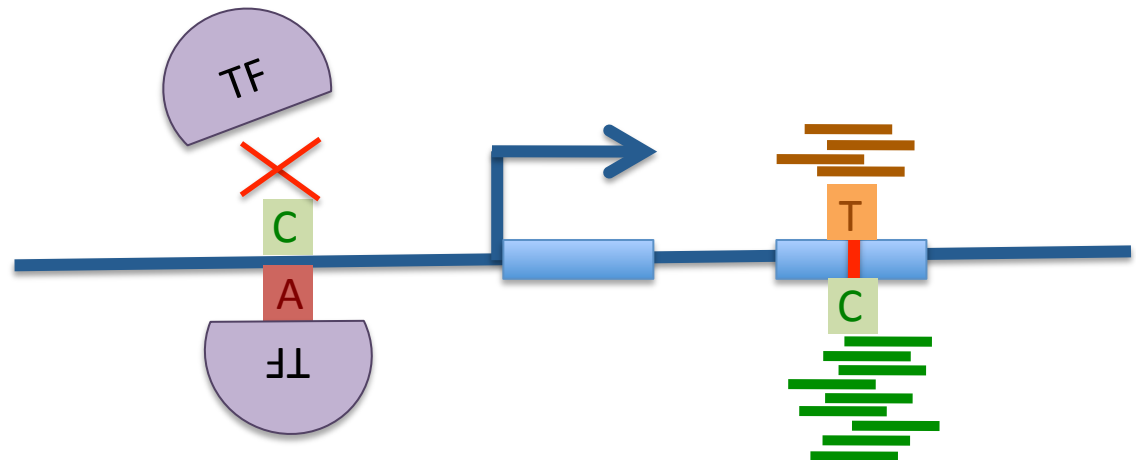
Find time-dependent allele specific expression to capture activated and inactivated genetic regulatory effects upon stimulation

# Condition specific allelic expression

Condition 1



Condition 2



# Assessing allele-specific expression in low input RNA-seq

## 1) Input data:

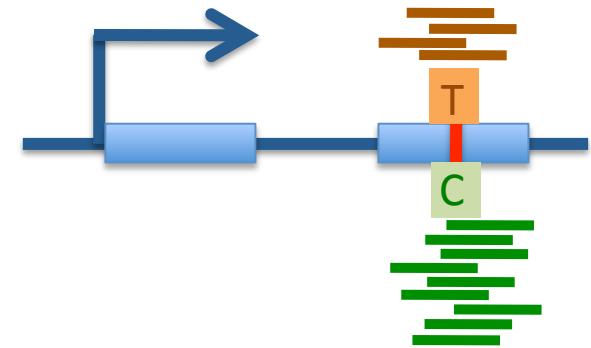
- Genotypes: called variants from RNA-seq
- RNA-seq uniquely aligned reads

## 2) Calculate read counts over heterozygous sites

- Pipeline from Kukurba *et al.* 2014

## 3) Filter and correct for sources of error

- Require BAS
- Remove mapping bias sites



$$\text{Ref ratio: } 4/(4+9) = 0.31$$

## 4) Analyze data

- Reference ratio : reference counts / total counts

# Assessing allele-specific expression in low input RNA-seq

## 1) Input data:

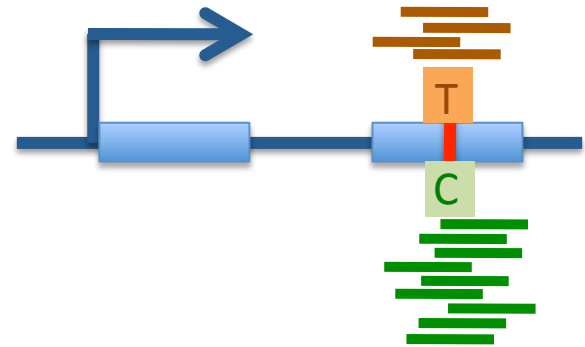
- Genotypes: called variants from RNA-seq
- RNA-seq uniquely aligned reads

☀ pipeline validation

## 2) Calculate read counts over heterozygous sites

- Pipeline from Kukurba *et al.* 2014

☀ weighted duplicate counts



## 3) Filter and correct for sources of error

- Require BAS
- Remove mapping bias sites

Ref ratio:  $4/(4+9) = 0.31$

☀ map to personalized masked genome

## 4) Analyze data

- Reference ratio : reference counts / total counts



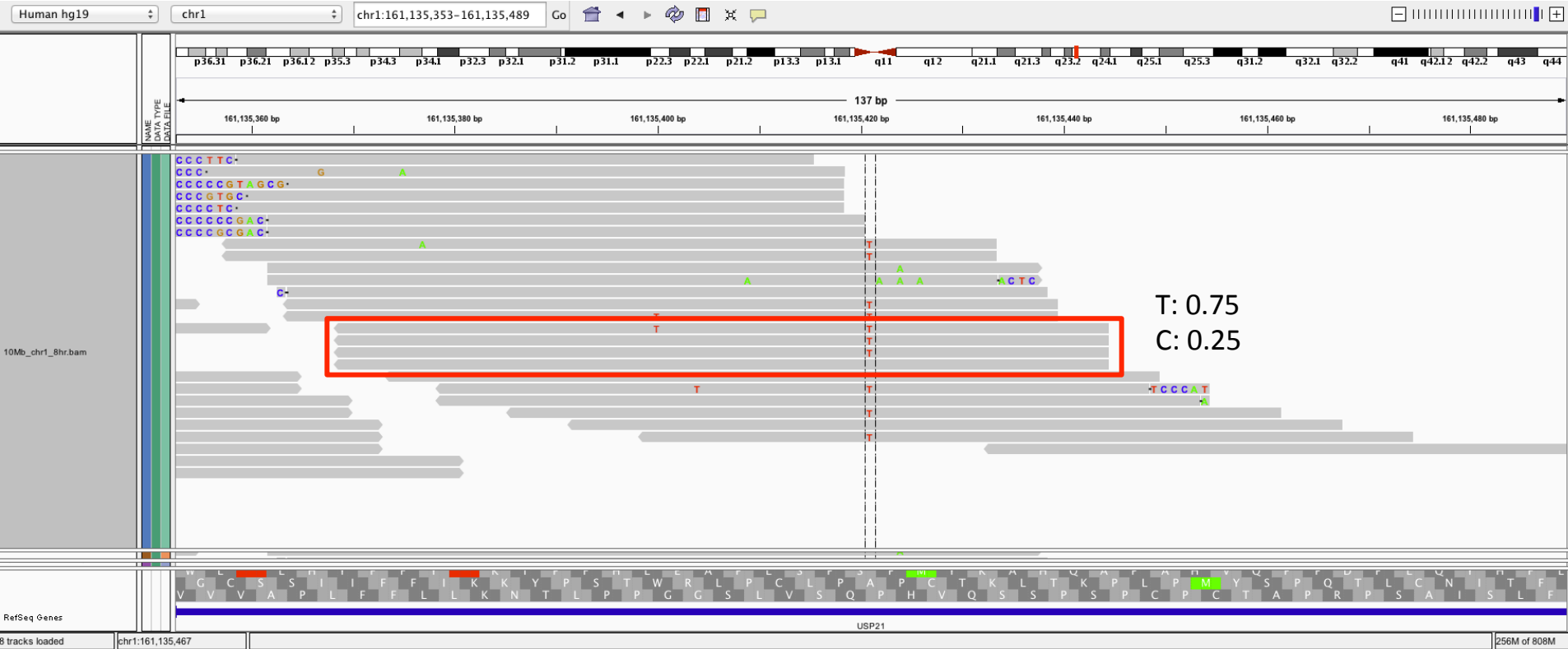
# Validating variant calling pipeline using NA12878

- “reference genotypes” are variants called by best practices of GATK from WGS
- mRNA-seq data from Kilpinen *et al.*, 2013: ~18M 50bp PE reads
- Differences with our RNA-seq datasets: less depth since only one condition assayed, paired-end data (so in reality ~36M reads), shorter reads

	GATK and Piskol filters	filts + overlap dbSNP	ASE tested*
<b>Total Hets Called</b>	13176	12566	4876
<b># in NA12878 calls</b>	12000	11976	4631
<b>with same alt allele</b>	11981	11957	4629
<b>with same genotype</b>	11911	11887	4624
<b>% concordant</b>	90.4	94.6	94.8
<b>% FP</b>	9.6	5.4	5.2

\*Both alleles seen + at least 10  
dup weighted counts per site

# Weighted duplicate counts



With dups

C: 7

T: 10

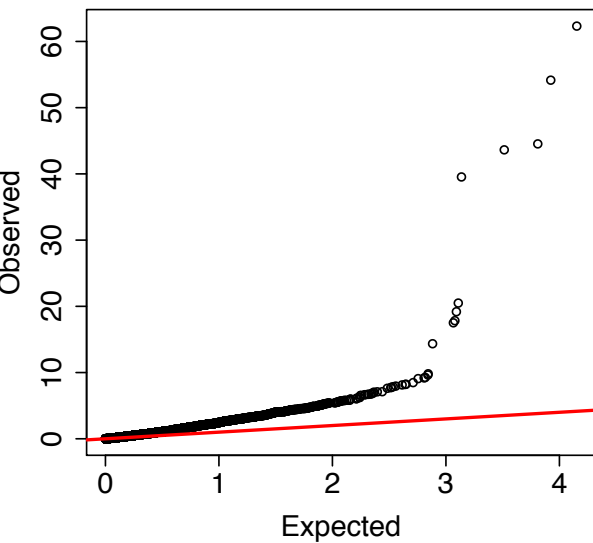
Weighted dups

C: 6.25

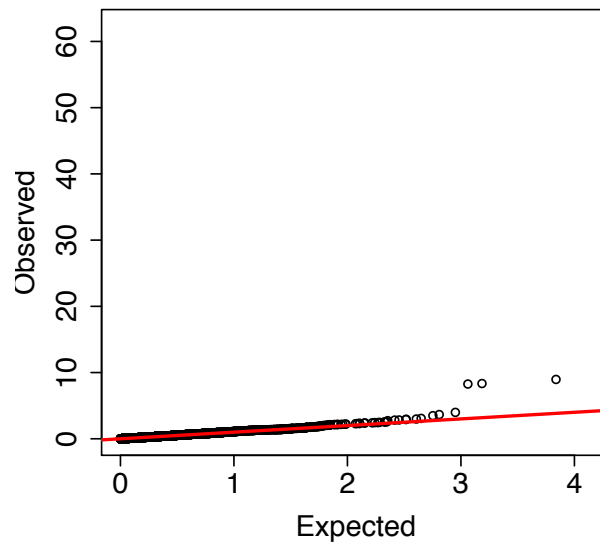
T: 5.75

# Logistic regression for time point vs allele counts

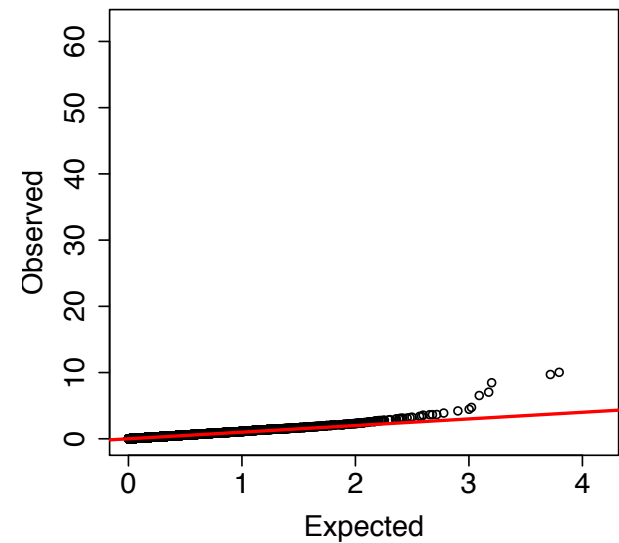
With duplicates  
 $\lambda = 3.1$



No duplicates  
 $\lambda = 1.1$

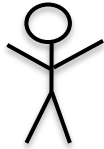


Weighted duplicates  
 $\lambda = 1.14$



\*time as quantitative variable

# Map to personalized masked genome



TACGCGATTCTGATCCGATAGC

TACGCGATTCTGATCCGATAGC

heterozygous site

reference genome

TACGCGATTCTGATCCGATAGC

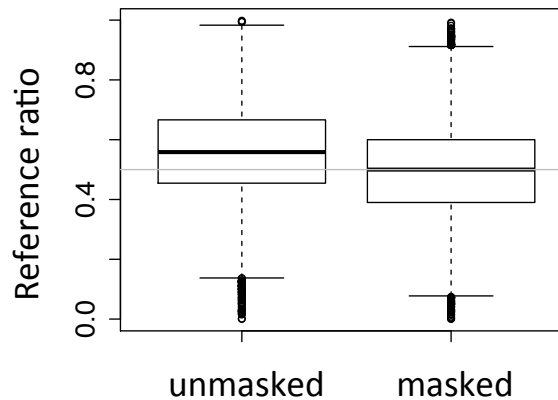


TACGCGATTCTNGATCCGATAGC

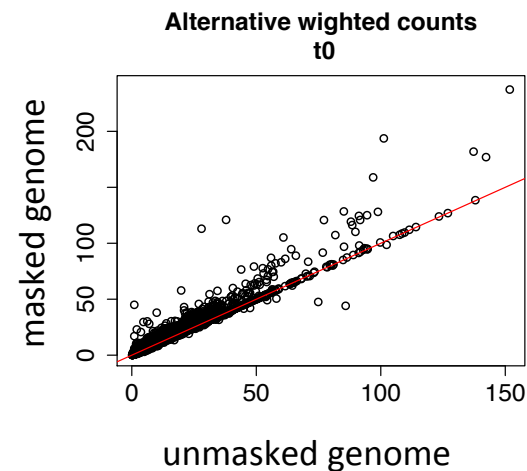
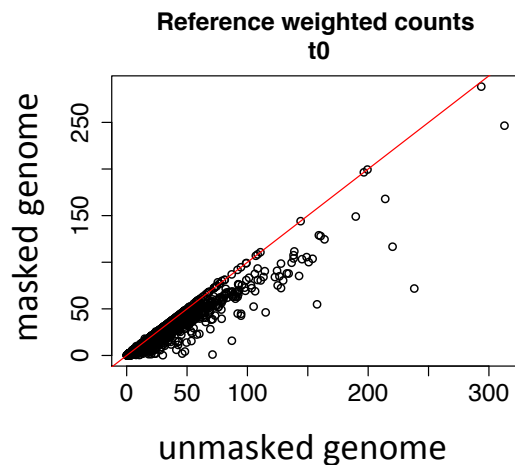
masked reference genome

# Map to personalized masked genome

Reference ratio  
distribution bias  
disappears



Reference allele tends to lose reads and alternative tends to gain reads

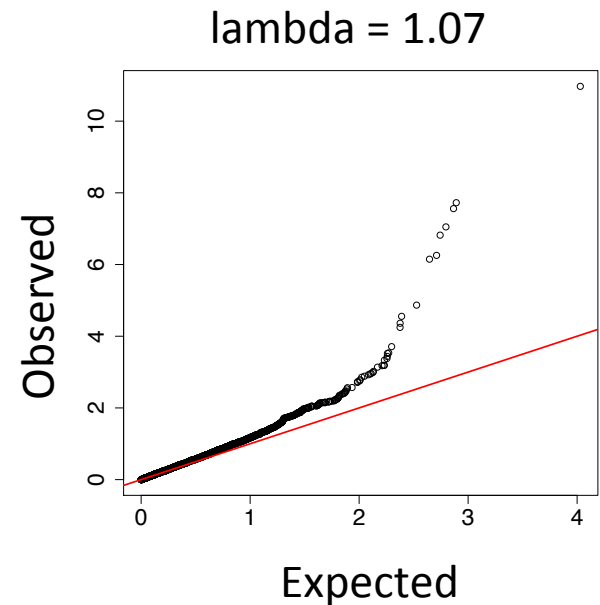
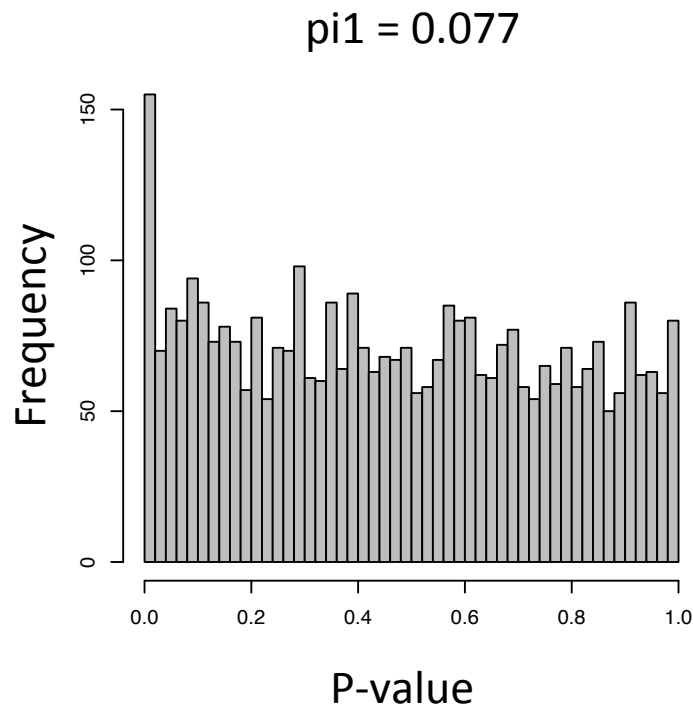


# Logistic regression time vs weighted allele counts

3,578 tested heterozygous sites

(BAS, min 10 weighted counts in all time points, masked genome)

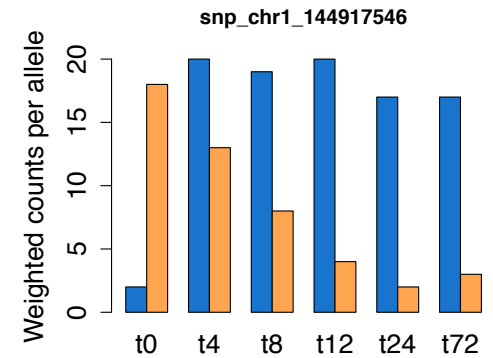
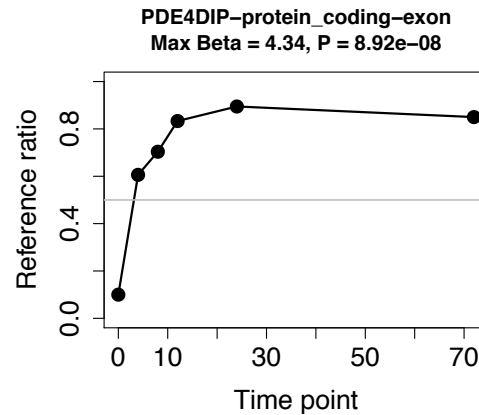
17 significant SNPs (FDR < 10%), spanning 13 protein coding genes, 1 non-coding site



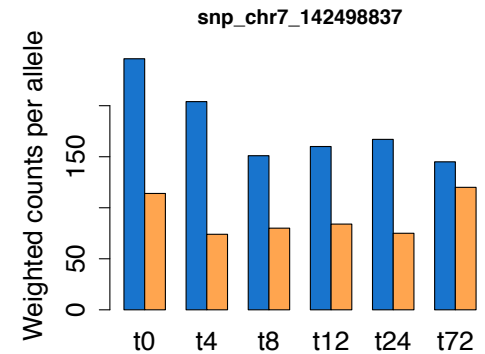
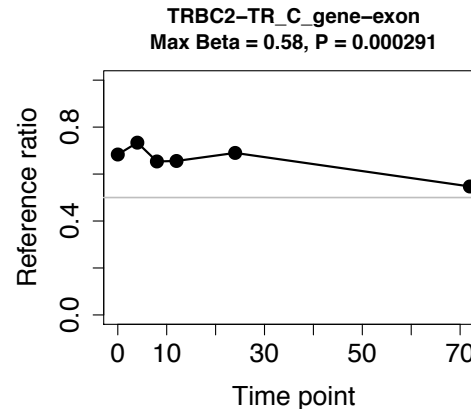
\*time as qualitative variable

# TASC ASE examples

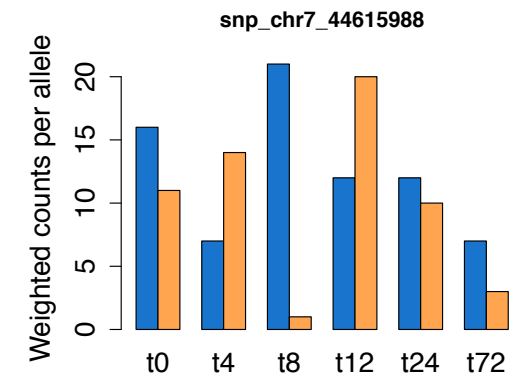
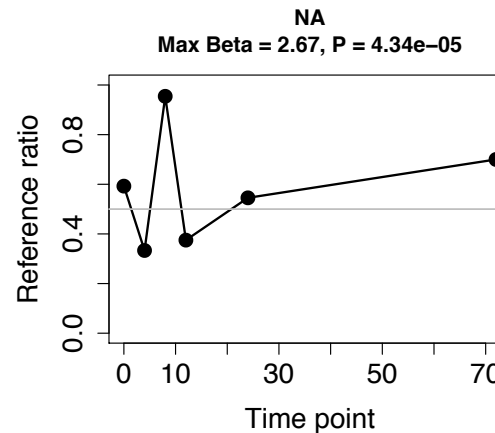
Anchors enzyme to Golgi  
Involved in eosinophilia



T cell receptor beta constant 2



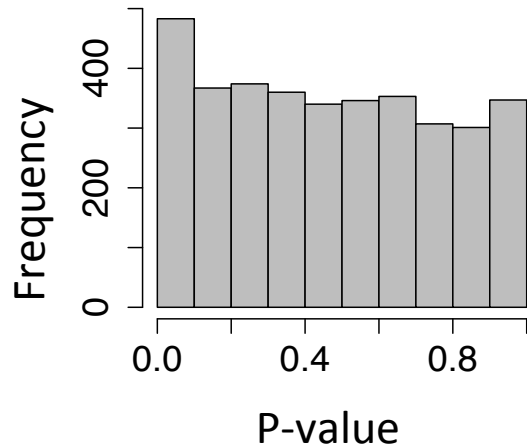
Non-coding  
Between 2 genes  
In Dnase I HS, 3 cell-types



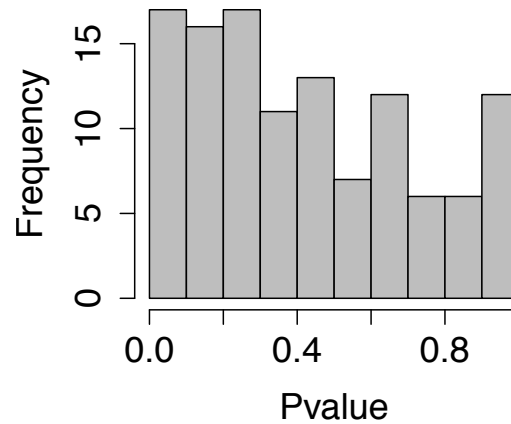
\*2 SNPs close to each other in 1 gene, and 1 case for chrX

# Enrichment of low P-values for RA genes

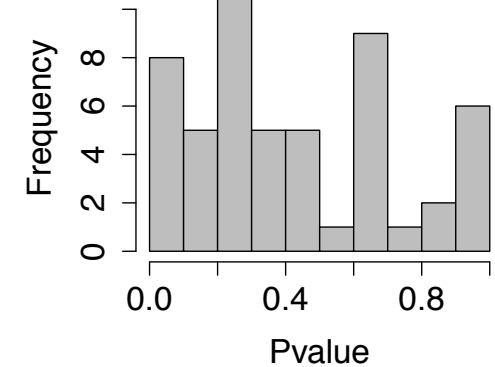
All tested sites  
 $\pi_1 = 0.077$   
 $N = 3578$



Hu *et al* genes with GWAS hit  
 $\pi_1 = 0.147$   
 $N = 117$

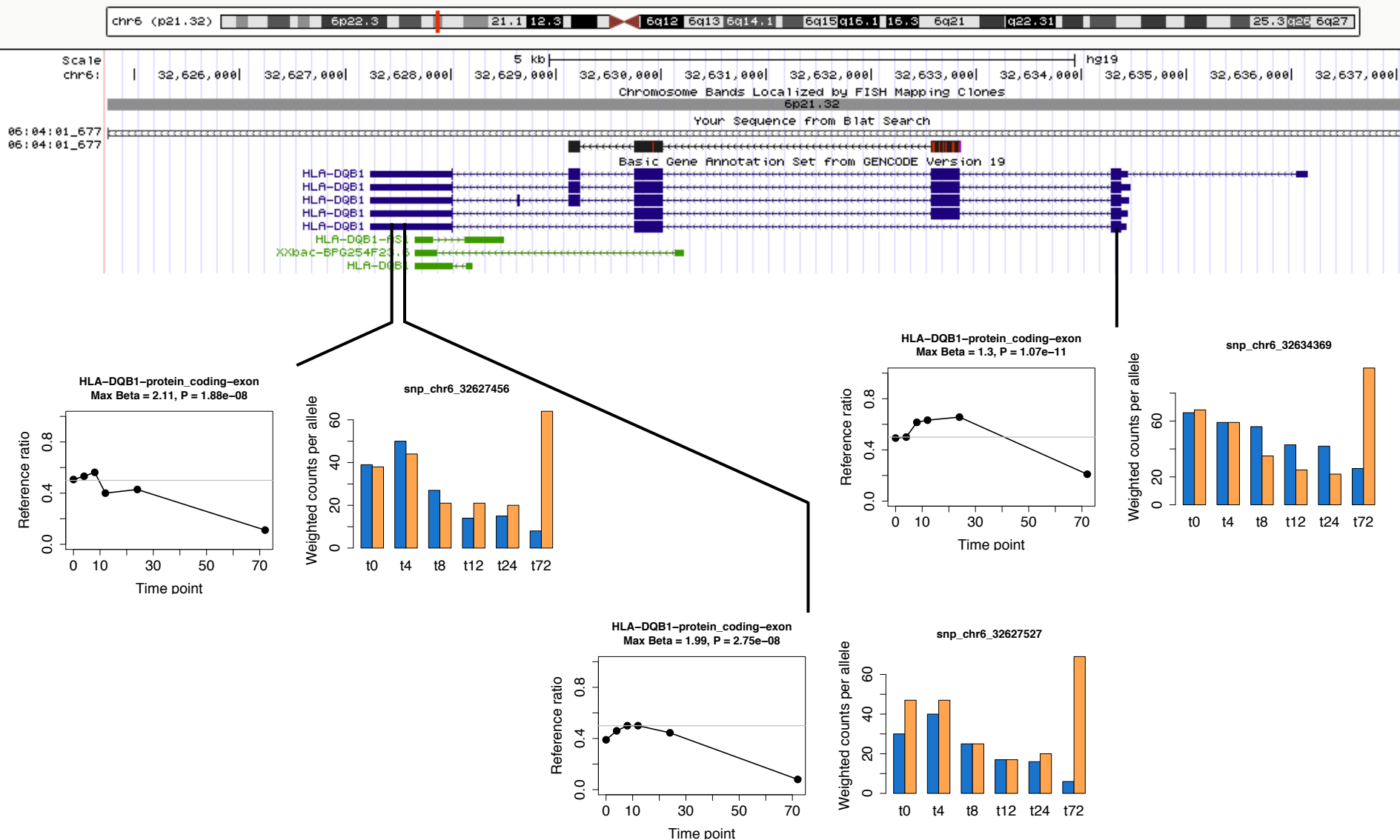


Vahedi *et al* super enhancer  
RA genes  
 $\pi_1 = 0.119$   
 $N = 153$





# HLA-DQB1: top TASC ASE signal



# HLA class II genotyping

## Clinical genotyping

Common alleles:

DQB1\*06:04

DQB1\*05:03

Rare alleles:

DQB1\*06:34/36/38/39/52/58/  
85/86/89/93/135/155/158N

DQB1\*05:06/08/10/13/15/16/23/  
24/38/39/40/41N/42/43/50/56/  
67/78

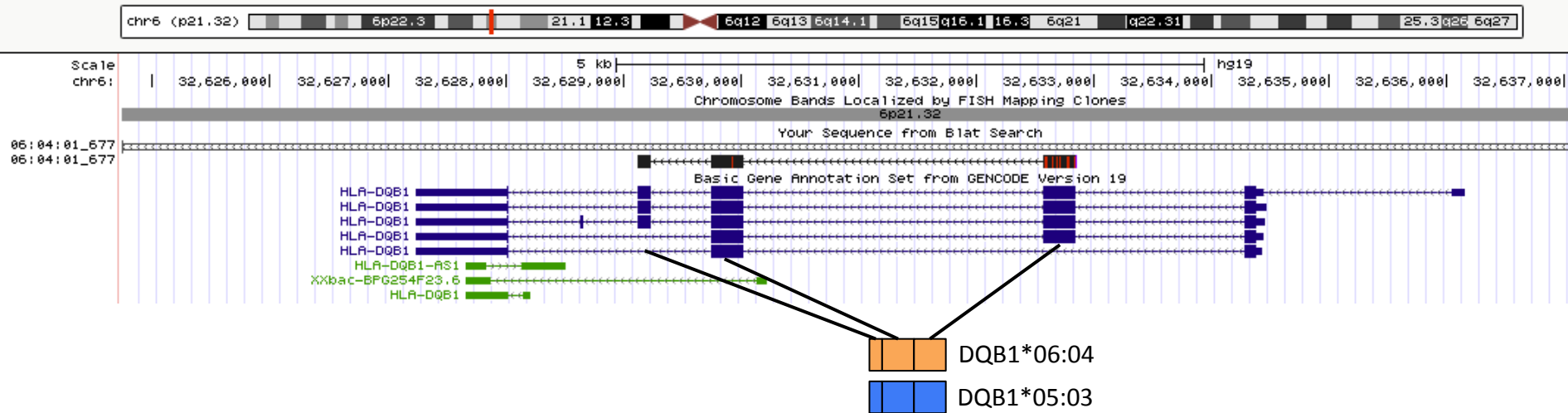
1) From each common allele, chose longest cDNA sequence available in ATHLATES IMGT/HLA db

2) Aligned these 2 sequences with BLAST and cut them so that they are the same length and well aligned (677bp, 96% id)

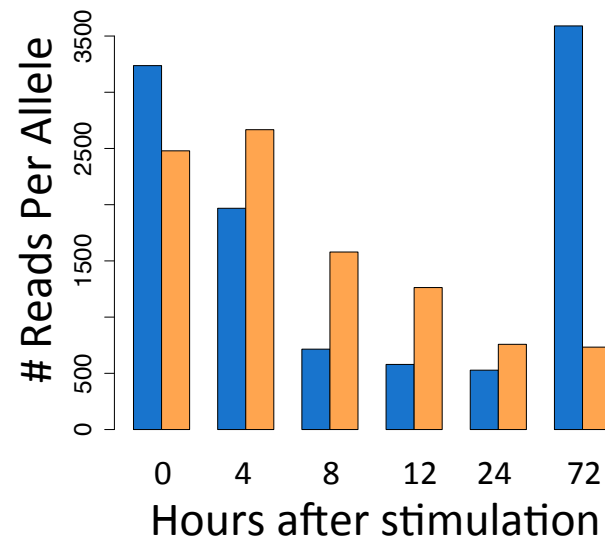
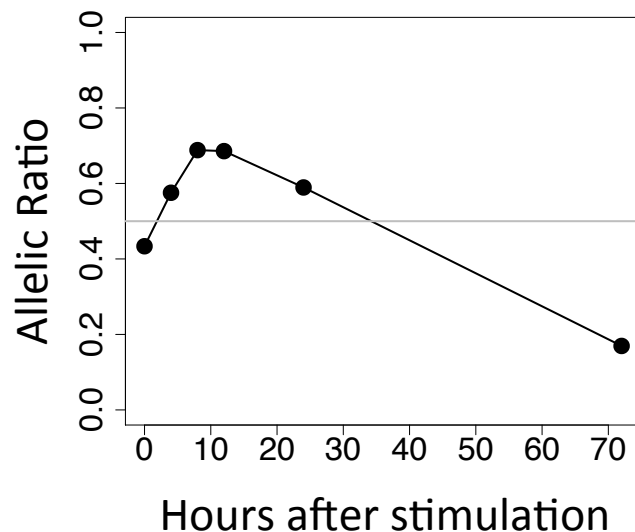
3) Aligned to ref genome with BLAT to get coordinates, and then masked these 3 exons

4) Added the 2 cDNA alleles to the masked ref genome and re-mapped all reads

# DQB1 TASC ASE replicates using 3 exons



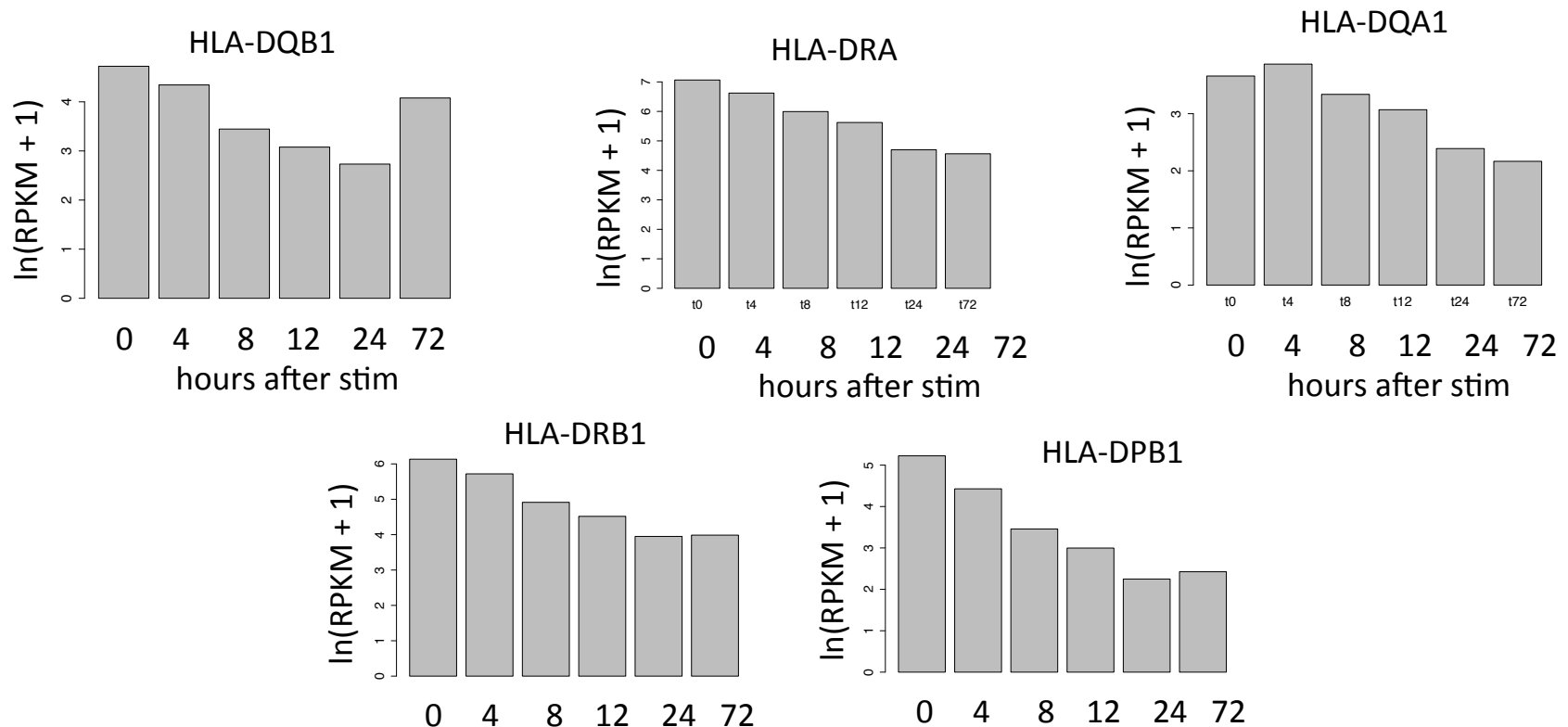
$P < 2.2e-16$



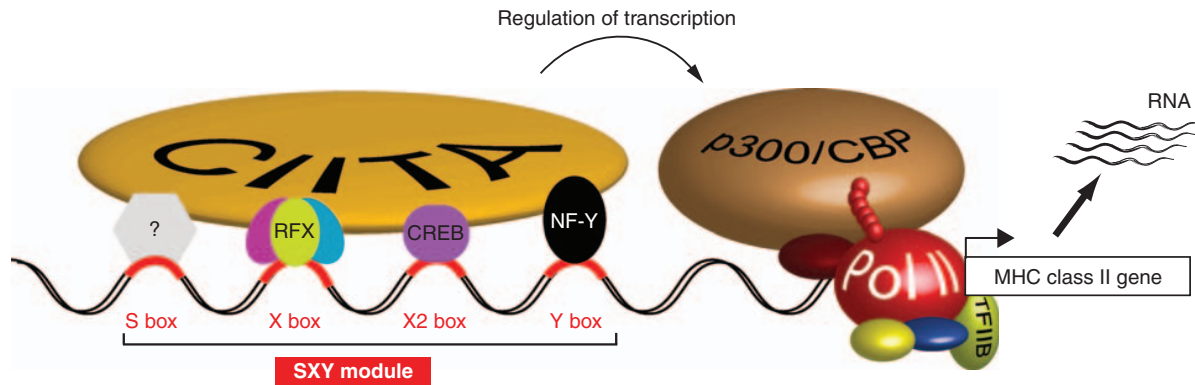
\*ASE for DQB1 has been reported in DCs, LCLs and melanoma cell line

# HLA class II genes in T cells

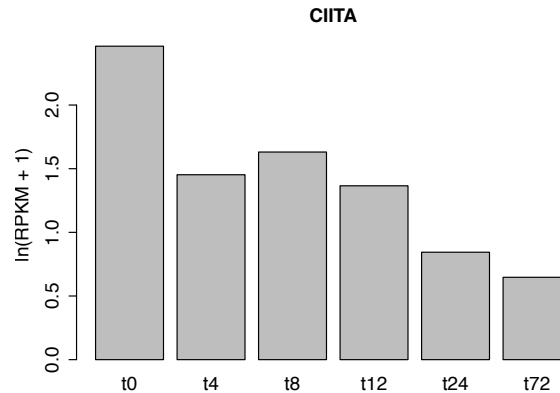
- HLA class II molecules are mainly expressed in professional antigen presenting cells (APCs)
- Normally not in T cells, but in activated human T cells yes (protein level)
- Expression can be induced in non-professional APCs
- Transcript expression values are counter-intuitive in our data:



# CIITA: co-activator of HLA class II genes

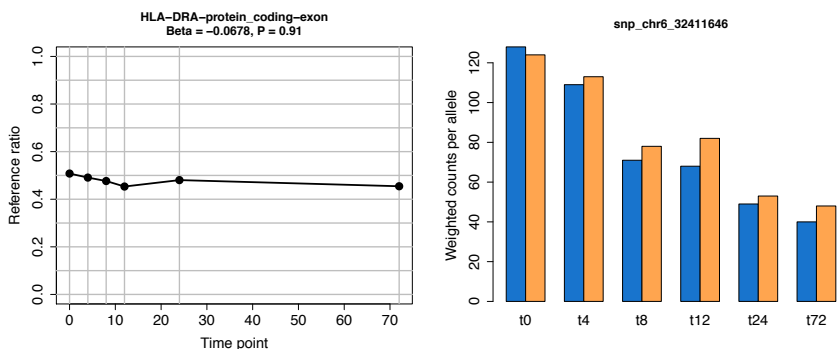


Handunnetthi *et al* 2010

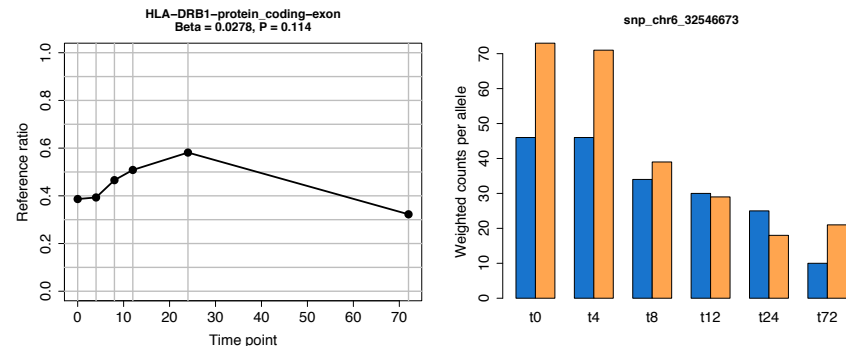


# No significant TASC ASE in other HLA class II genes

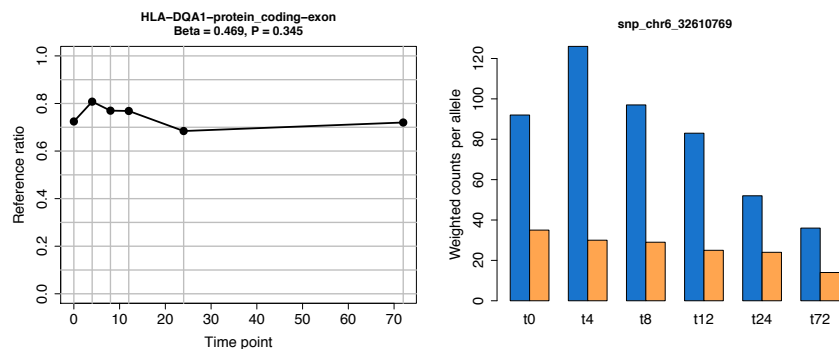
## HLA-DRA



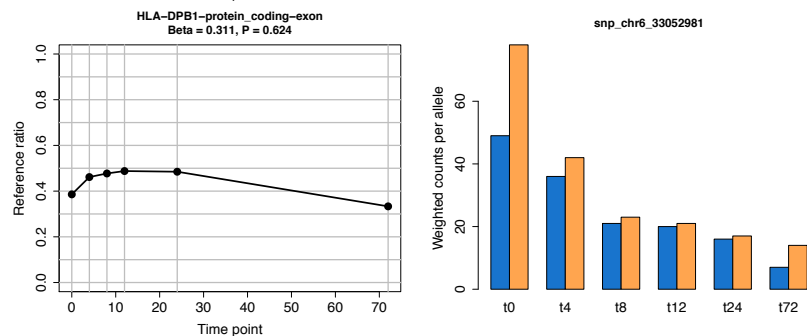
## HLA-DRB1



## HLA-DQA1



## HLA-DPB1



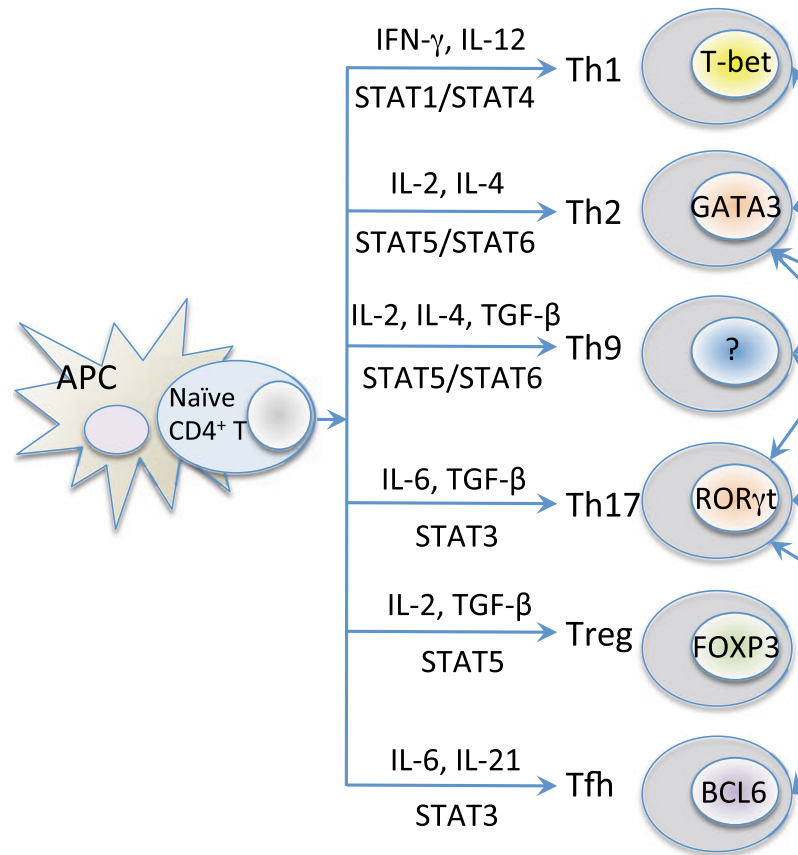
# Next steps for time-dependent ASE and DQB1

- One-off experiment to check if protein expr DQB1 goes up upon CD4+ Tmem stimulation (can we at the same time identify which cellular subtype the signal is coming from?)
- Check transcript levels for HLA class II genes in other T cell datasets (will check in Deepak's)
- Make sure DQB1 signal does not come from DQB2 (unlikely)
- Measure TASC ASE in other genes by phasing SNPs and fusing exons if possible (more power for log reg)
- Expand study to many genotyped individuals (chose based on DQB1 genotypes?, same time points or less?)

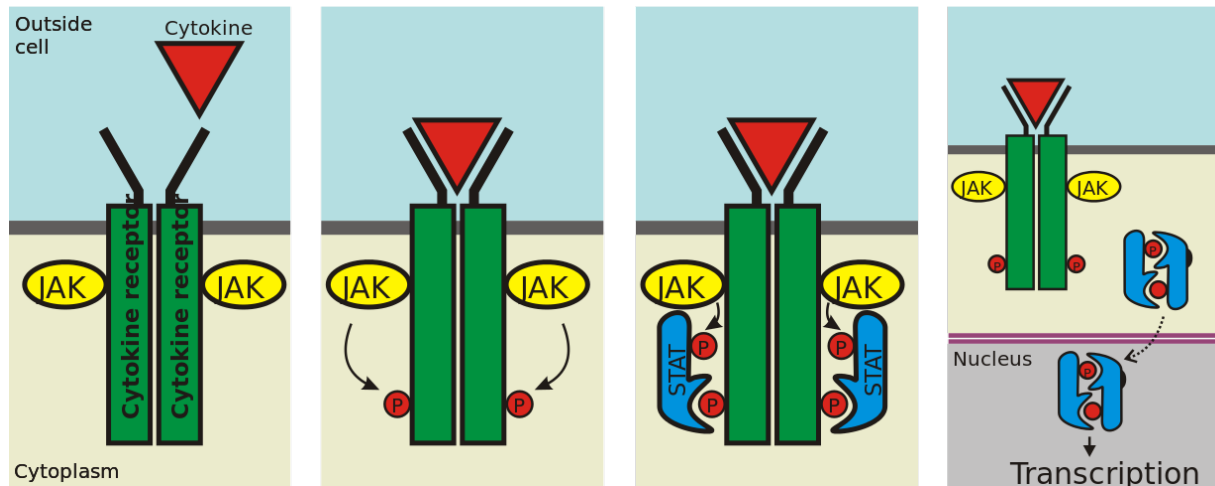
# STAT PILOT

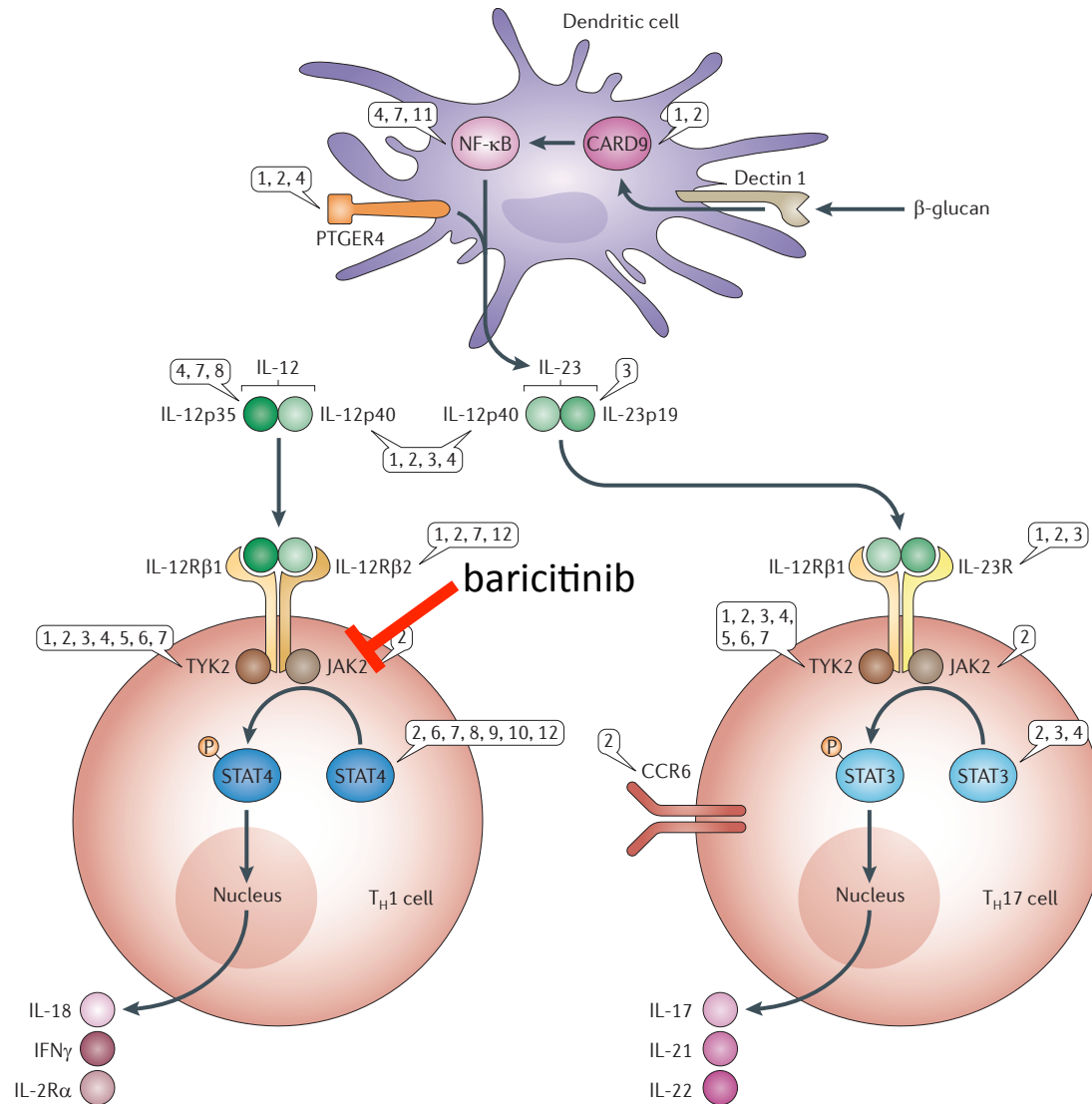


# Signal Transducer and Activator of Transcription (STATs)



# STAT TFs are activated by Janus Kinase (JAK)

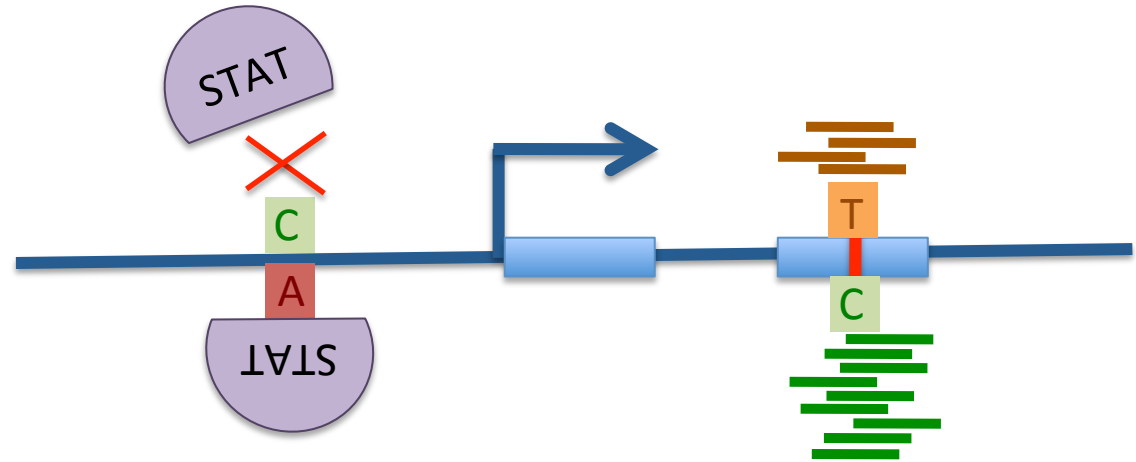




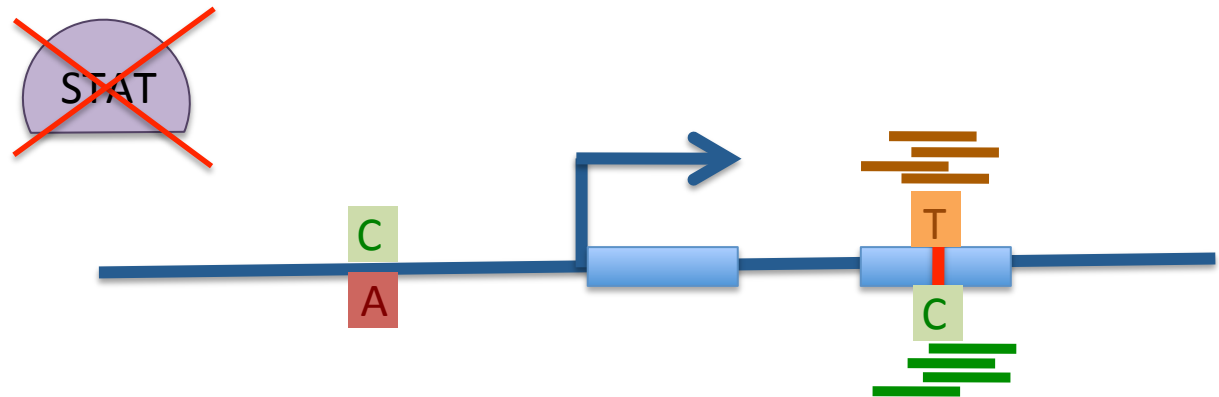
- |                              |                             |                                |
|------------------------------|-----------------------------|--------------------------------|
| ① Ankylosing spondylitis     | ⑤ Type 1 diabetes           | ⑨ Systemic lupus erythematosus |
| ② Inflammatory bowel disease | ⑥ Rheumatoid arthritis      | ⑩ Sjögren's syndrome           |
| ③ Psoriasis                  | ⑦ Primary biliary cirrhosis | ⑪ Ulcerative colitis           |
| ④ Multiple sclerosis         | ⑧ Coeliac disease           | ⑫ Behçet's disease             |

# Looking for genetic differential regulation mediated by STATs

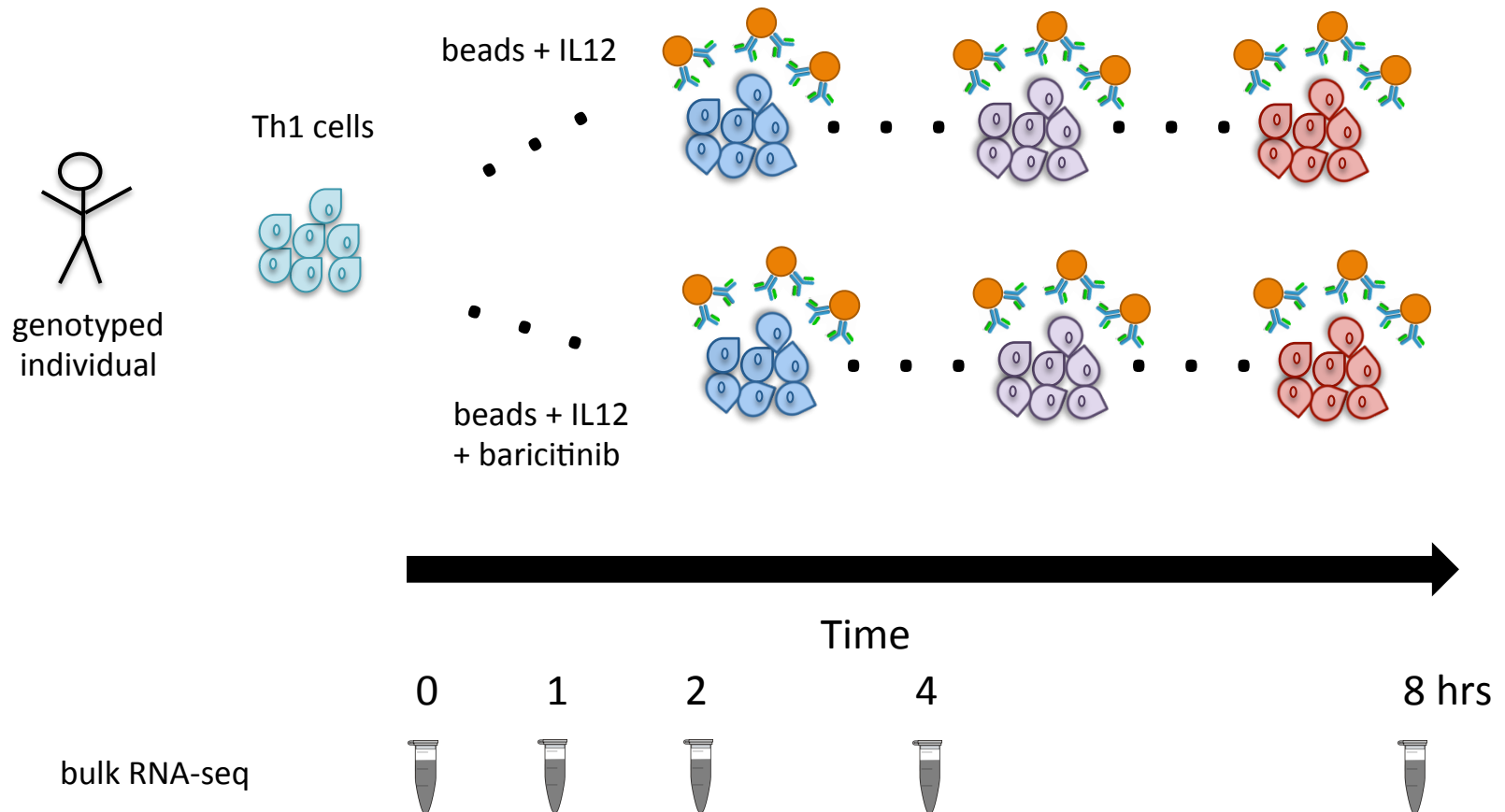
Condition 1



Condition 2



# Pilot experiment design



# Pre-experiments

- RT-PCR

At which time points transcription levels of direct STAT4 target genes peak?

- IFNgama

- TNF

- IL12RB2/Tbet

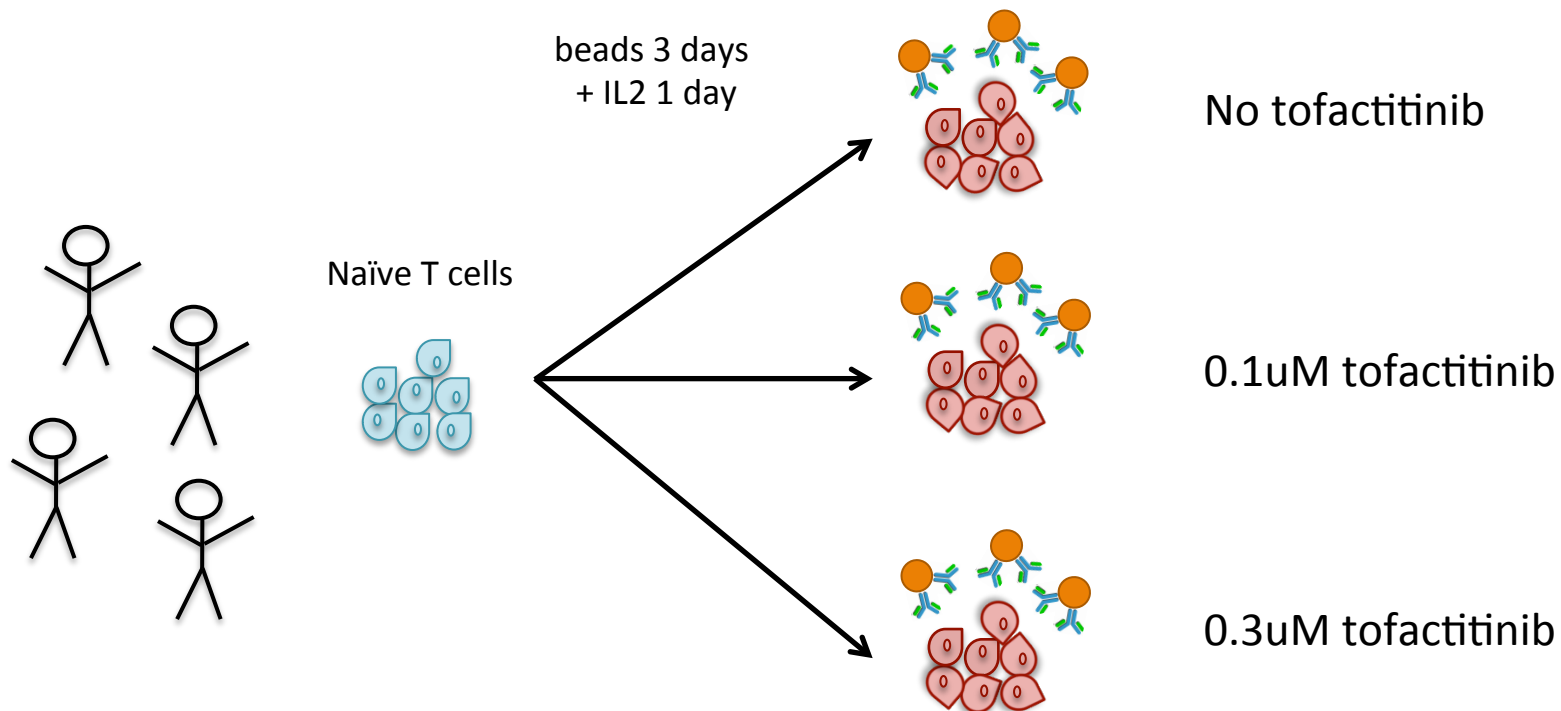
- Phospho flow

Which STATs get phosphorylated in Th1 upon stimulation

Which STATs are affected by baricitinib?

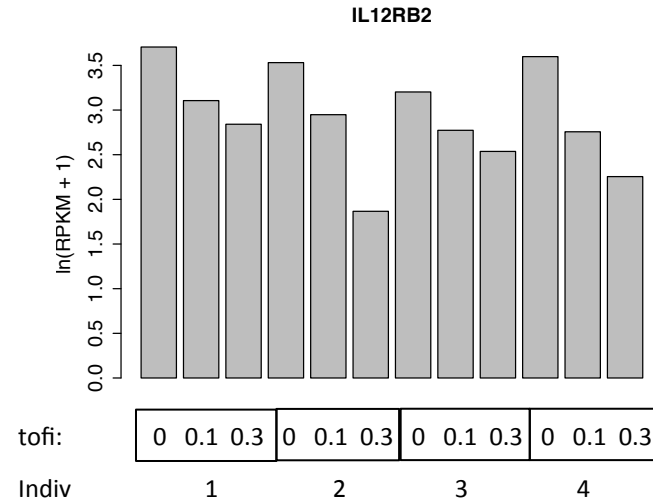
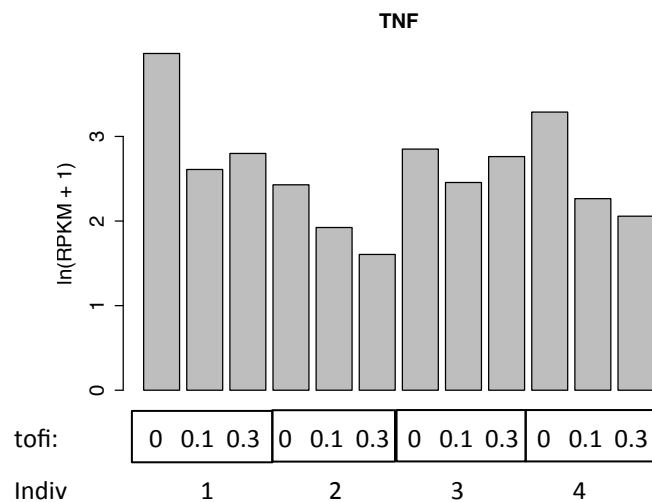
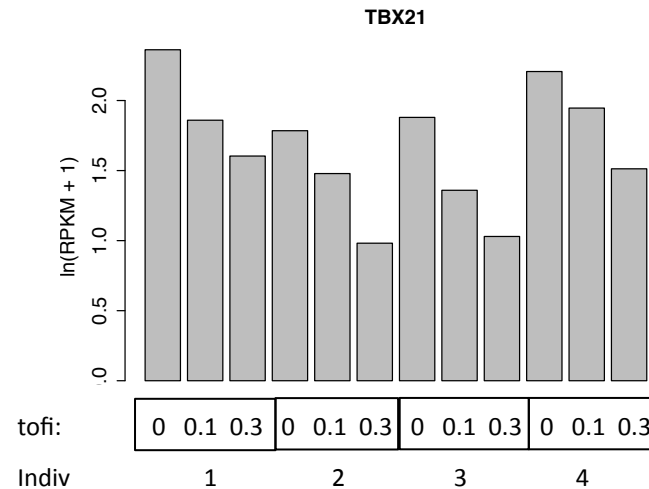
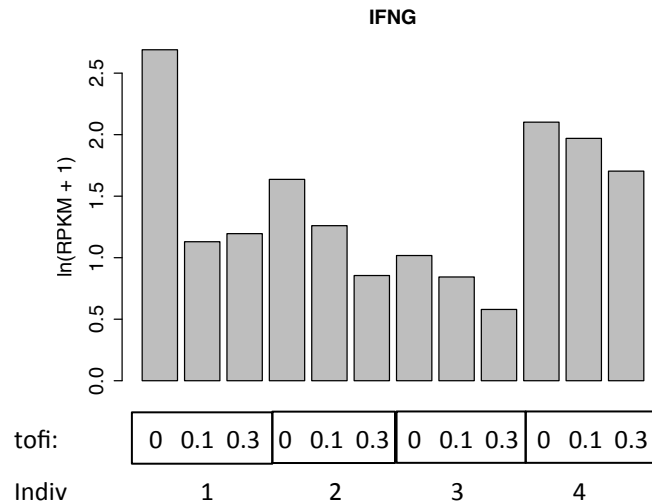
# Analyses on published dataset using tofi on T cells

Vahedi *et al* 2015



RNA-seq at 3 days, for 3 conditions and 4 indivs

# STAT4 target genes are mostly down-regulated by tofi

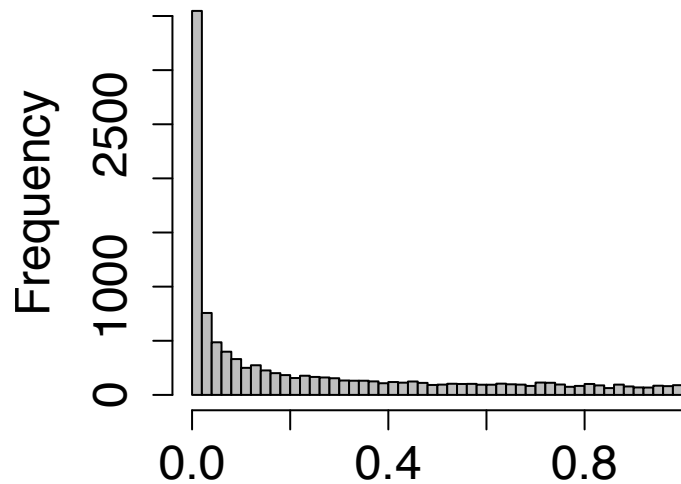




# Tofi regulated genes are enriched for RA genes

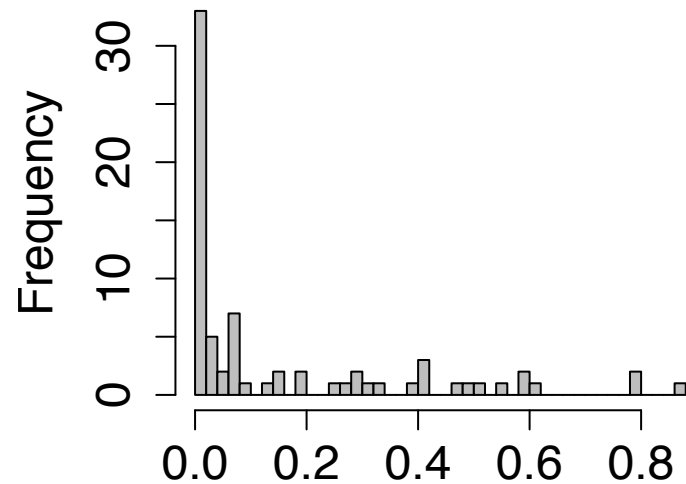
P-values for differential expression analysis

All 11K expressed genes  
 $\pi_1 = 0.645$



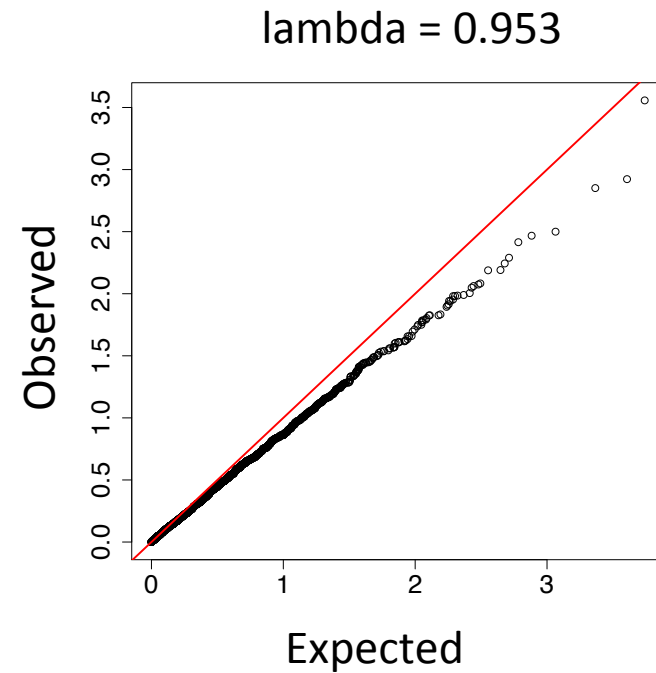
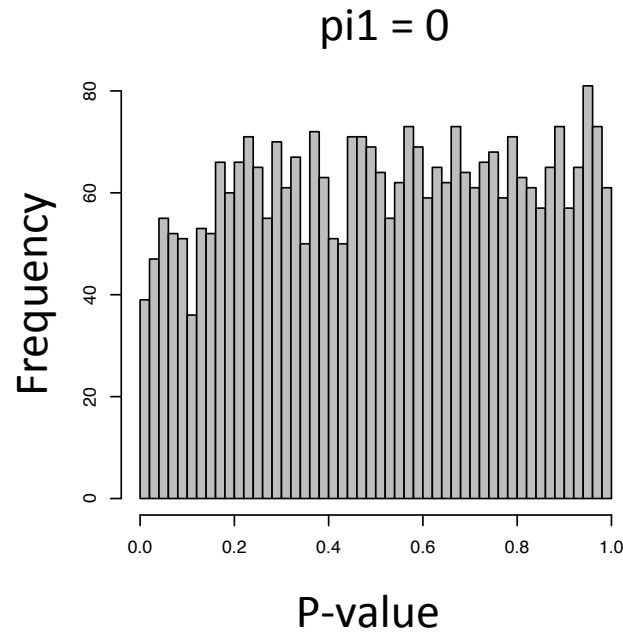
Diff Expr (Tofi vs noTofi) P-value

RA genes Okada *et al*  
 $\pi_1 = 0.966$



Diff Expr (Tofi vs noTofi) P-value

# No enrichment for significant Tofi-dep ASE events



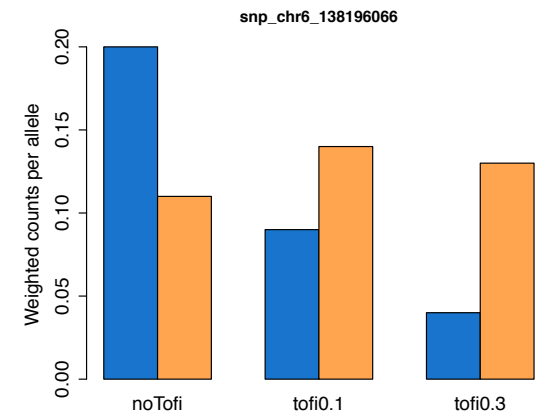
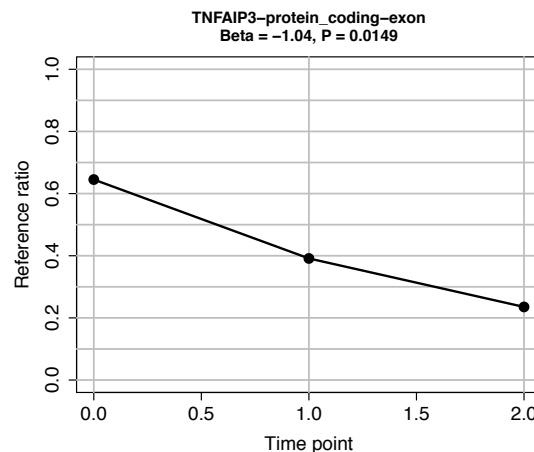
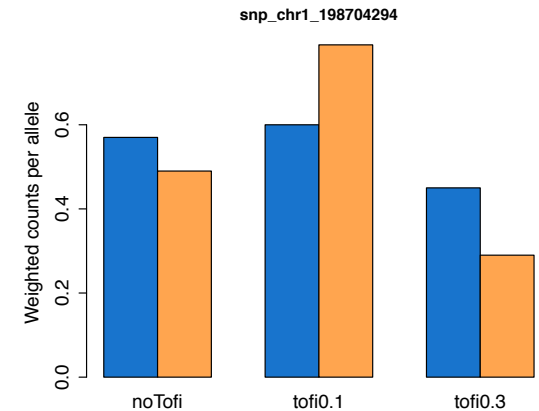
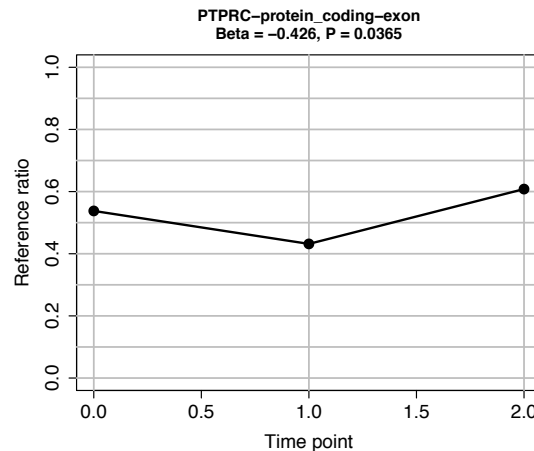
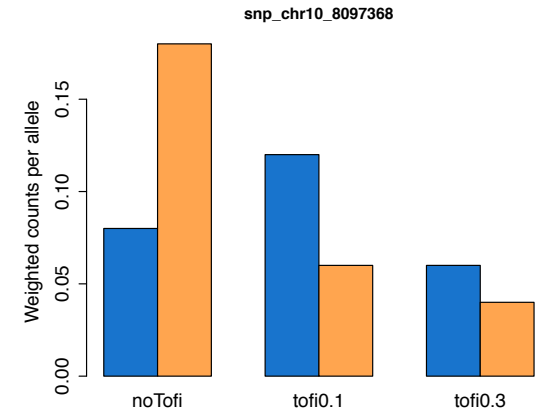
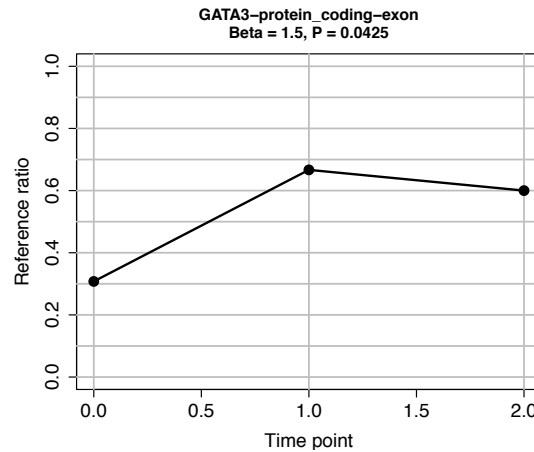
N = 3,090 tested sites

# 3 SNPs in RA genes with $P < 0.05$

GATA3 TF

PTPRC  
(CD45, regulator of  
cytokine receptor  
signaling)

TNF induced protein  
(zinc finger, inhibits NFKB  
activation)



# Conclusions

Condition specific ASE can be detected in a single individual (if sequenced at enough depth and at enough conditions)

Condition specific ASE can be used to detect genetic regulatory effects of disease genes and dissect at which cellular states they are active

Condition specific ASE in DQB1 seems to be real (but needs replication in additional individuals)