

polyfid

Load required libraries, load and merge MAGERI results.

```
library(plyr); library(ggplot2); library(seqLogo); library(Biostrings); library(reshape2); library(gplots)
```

```
## Warning: package 'ggplot2' was built under R version 3.2.4
```

```
## Loading required package: grid
```

```
## Loading required package: BiocGenerics
```

```
## Loading required package: parallel
```

```
##
```

```
## Attaching package: 'BiocGenerics'
```

```
## The following objects are masked from 'package:parallel':
```

```
##
```

```
##   clusterApply, clusterApplyLB, clusterCall, clusterEvalQ,  
##   clusterExport, clusterMap, parApply, parCapply, parLapply,  
##   parLapplyLB, parRapply, parSapply, parSapplyLB
```

```
## The following objects are masked from 'package:stats':
```

```
##
```

```
##   IQR, mad, xtabs
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
##   anyDuplicated, append, as.data.frame, as.vector, cbind,  
##   colnames, do.call, duplicated, eval, evalq, Filter, Find, get,  
##   grep, grepl, intersect, is.unsorted, lapply, lengths, Map,  
##   mapply, match, mget, order, paste, pmax, pmax.int, pmin,  
##   pmin.int, Position, rank, rbind, Reduce, rownames, sapply,  
##   setdiff, sort, table, tapply, union, unique, unlist, unsplit
```

```
## Loading required package: S4Vectors
```

```
## Loading required package: stats4
```

```
##
```

```
## Attaching package: 'S4Vectors'
```

```
## The following object is masked from 'package:plyr':
```

```
##
```

```
##   rename
```

```
## Loading required package: IRanges
```

```
##
## Attaching package: 'IRanges'

## The following object is masked from 'package:plyr':
##
##      desc

## Loading required package: XVector

##
## Attaching package: 'XVector'

## The following object is masked from 'package:plyr':
##
##      compact

## Warning: package 'gplots' was built under R version 3.2.4

##
## Attaching package: 'gplots'

## The following object is masked from 'package:IRanges':
##
##      space

## The following object is masked from 'package:stats':
##
##      lowess
```

```
df.meta <- read.table("metadata.txt", header=T, sep = "\t")
```

```
df <- data.frame()
```

```
for (project in levels(df.meta$project)) {
  for (sample in levels(df.meta$sample)) {
    fname <- paste(project, sample, "variant.caller.txt", sep = ".")
    if (file.exists(fname)) {
      df.1 <- read.table(fname, header=T, sep="\t")
      df.1$project <- project
      df.1$sample <- sample
      df <- rbind(df, df.1)
    }
  }
}
```

```
df$project <- as.factor(df$project)
df$sample <- as.factor(df$sample)
```

```
df <- merge(df, df.meta, all.x=T, all.y=F)
```

```
template <- "TGGGATCCATTATCGGCGGCGAATTTACCACCATTGAAAACAGCCGTGGTTTGGCGCGATTATCGTCGTCATCGTGGCGGCAGCGTGA
```

Error rates

Overall error rates

```
df.er <- ddpoly(df, .(project, name), summarize,
  mismatches = sum(count.major), umi.count = round(mean(coverage)),
  err.rate = mismatches / umi.count / nchar(template) / mean(cycles),
  delta = sqrt(mismatches / umi.count * (1 - mismatches / umi.count) / umi.count) / nchar(
  err.lb = err.rate - 1.96 * delta, err.ub = err.rate + 1.96 * delta)

print(df.er)
```

##	project	name	mismatches	umi.count	err.rate
## 1	polerr73	Encyclo	7759	558336	5.428379e-06
## 2	polerr73	Kappa HF	2444	240548	3.968802e-06
## 3	polerr73	Phusion	2902	512908	1.768105e-06
## 4	polerr73	SD-HS	1696	72692	9.113795e-06
## 5	polerr73	SNP-detect	303	57755	2.049336e-06
## 6	polerr73	Taq-HS	1351	78646	6.710251e-06
## 7	polerr73	Tersus	2430	180029	5.272588e-06
## 8	polerr73	Tersus-SNP-buffer	2035	488168	1.628378e-06
## 9	polerr73	TruSeq	211	43380	1.899997e-06
## 10	polerr73	Velox	6089	427259	5.566918e-06
## 11	polerr82	Encyclo	3127	200300	6.098274e-06
## 12	polerr82	Kappa HF	1040	116021	3.501521e-06
## 13	polerr82	Phusion	3349	471735	2.218539e-06
## 14	polerr82	SD-HS	2294	94810	9.451469e-06
## 15	polerr82	SNP-detect	372	92765	1.566458e-06
## 16	polerr82	Taq-HS	3076	133952	8.970097e-06
## 17	polerr82	Tersus	1655	497406	1.299712e-06
## 18	polerr82	Tersus-SNP-buffer	1391	303834	1.788343e-06
## 19	polerr82	TruSeq	483	113803	1.657881e-06
## 20	polerr82	Velox	2234	141485	6.167836e-06
##	delta	err.lb	err.ub		
## 1	6.119677e-08	5.308433e-06	5.548325e-06		
## 2	7.987141e-08	3.812255e-06	4.125350e-06		
## 3	3.272858e-08	1.703957e-06	1.832253e-06		
## 4	2.187056e-07	8.685132e-06	9.542458e-06		
## 5	1.174220e-07	1.819188e-06	2.279483e-06		
## 6	1.809874e-07	6.355515e-06	7.064986e-06		
## 7	1.062355e-07	5.064366e-06	5.480810e-06		
## 8	3.602183e-08	1.557775e-06	1.698980e-06		
## 9	1.304827e-07	1.644251e-06	2.155743e-06		
## 10	7.083125e-08	5.428088e-06	5.705747e-06		
## 11	1.081998e-07	5.886203e-06	6.310346e-06		
## 12	1.080898e-07	3.289665e-06	3.713377e-06		
## 13	3.819992e-08	2.143667e-06	2.293411e-06		
## 14	1.949324e-07	9.069401e-06	9.833536e-06		
## 15	8.105407e-08	1.407592e-06	1.725324e-06		
## 16	1.598672e-07	8.656758e-06	9.283437e-06		
## 17	3.189512e-08	1.237197e-06	1.362226e-06		
## 18	4.783996e-08	1.694577e-06	1.882109e-06		
## 19	7.527596e-08	1.510340e-06	1.805422e-06		
## 20	1.294599e-07	5.914094e-06	6.421577e-06		

Error rate consistency

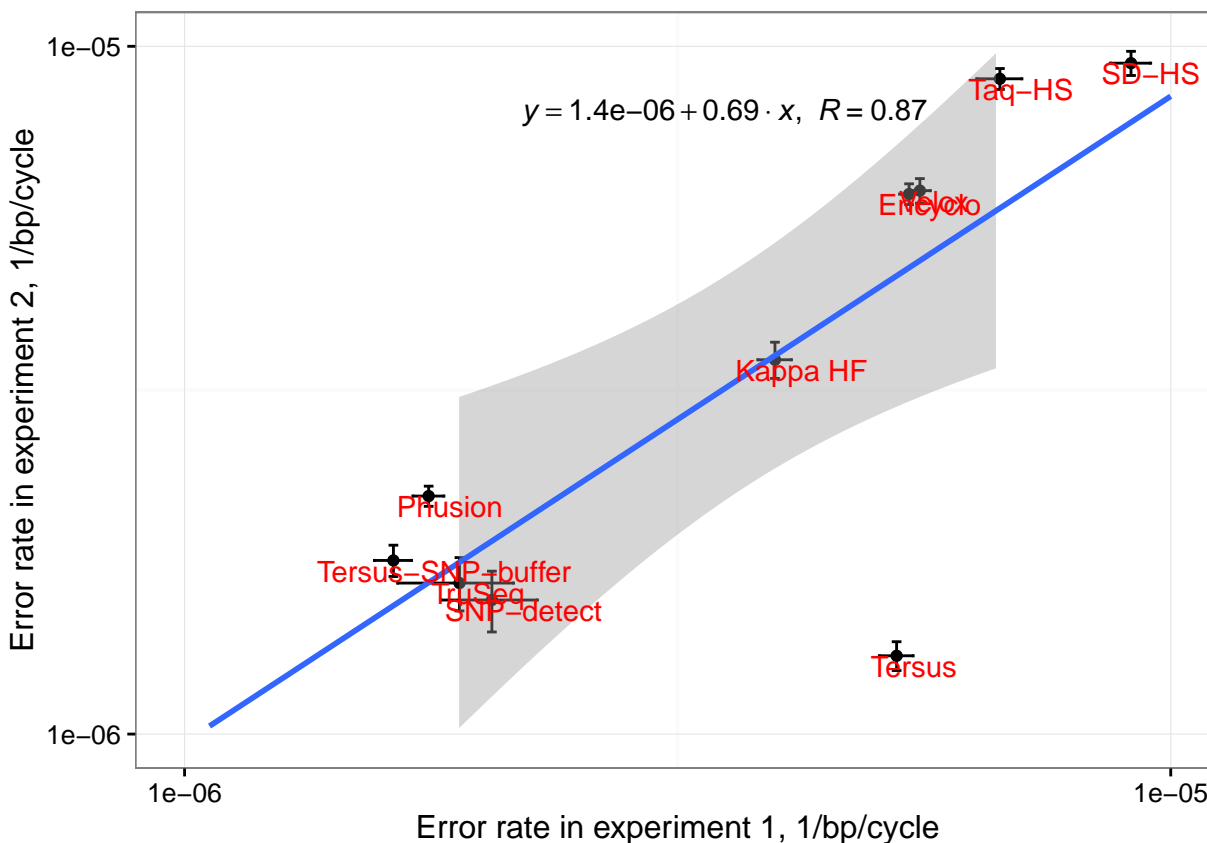
```
df.er.cast <- dcast(df.er, name ~ project,
                    value.var = "err.rate")
df.er.cast2 <- dcast(df.er, name ~ project,
                     value.var = "delta")

df.er.cast <- merge(df.er.cast, df.er.cast2, by = "name")

m <- lm(polerr73.x ~ polerr82.x, df.er.cast);
eq <- substitute(italic(y) == a + b %.% italic(x)*", "~~italic(R)~"~r,
                 list(a = format(coef(m)[1], digits = 2),
                       b = format(coef(m)[2], digits = 2),
                       r = format(sqrt(summary(m)$r.squared), digits = 2)))
lbl<-as.character(as.expression(eq))

ggplot(df.er.cast, aes(x=polerr73.x, y=polerr82.x)) + geom_point() +
  geom_errorbarh(aes(xmax=polerr73.x+1.96*polerr73.y,xmin=polerr73.x-1.96*polerr73.y)) +
  geom_errorbar(aes(ymax=polerr82.x+1.96*polerr82.y,ymin=polerr82.x-1.96*polerr82.y)) +
  geom_smooth(method = "lm", fullrange = T) +
  geom_text(aes(label=name), vjust=1, hjust = .3, color="red") +
  annotate("text", x = 2e-6, y = 8e-6, label = lbl, hjust=-0.1, parse = TRUE)+
  scale_x_log10(name="Error rate in experiment 1, 1/bp/cycle", limits=c(1e-6,1e-5)) +
  scale_y_log10(name="Error rate in experiment 2, 1/bp/cycle", limits=c(1e-6,1e-5)) +
  theme_bw()
```

```
## Warning: Removed 2 rows containing missing values (geom_smooth).
```



Combine error rates from a pair of experiments.

```
df.coverage.summary <- ddply(df, .(name, project), summarize, coverage = mean(coverage))
df.coverage.summary <- ddply(df.coverage.summary, .(name), summarize, coverage = sum(coverage))

df <- ddply(df, .(name, mutation), summarize, count.major = sum(count.major))

df <- merge(df, df.coverage.summary, by = "name")
```

Error substitution patterns

Parse mutation signatures (needed for further analysis)

```
df$mut.split <- sapply(df$mutation, function(x) strsplit(as.character(x), "[S:>]"))
df$mutation.pos <- as.integer(sapply(df$mut.split, function(x) x[2]))
df$mutation.from <- sapply(df$mut.split, function(x) x[3])
df$mutation.to <- sapply(df$mut.split, function(x) x[4])
df$mut.split <- NULL
```

Substitution signature preferences

```
df$mutation.signature <- paste(df$mutation.from, df$mutation.to, sep = ">")

sign.rep <- data.frame(mutation.signature = c("A>C", "A>G", "A>T", "C>A", "C>G", "C>T", "G>A", "G>C", "G>T", "T>A", "T>C", "T>G", "T>T"))
```

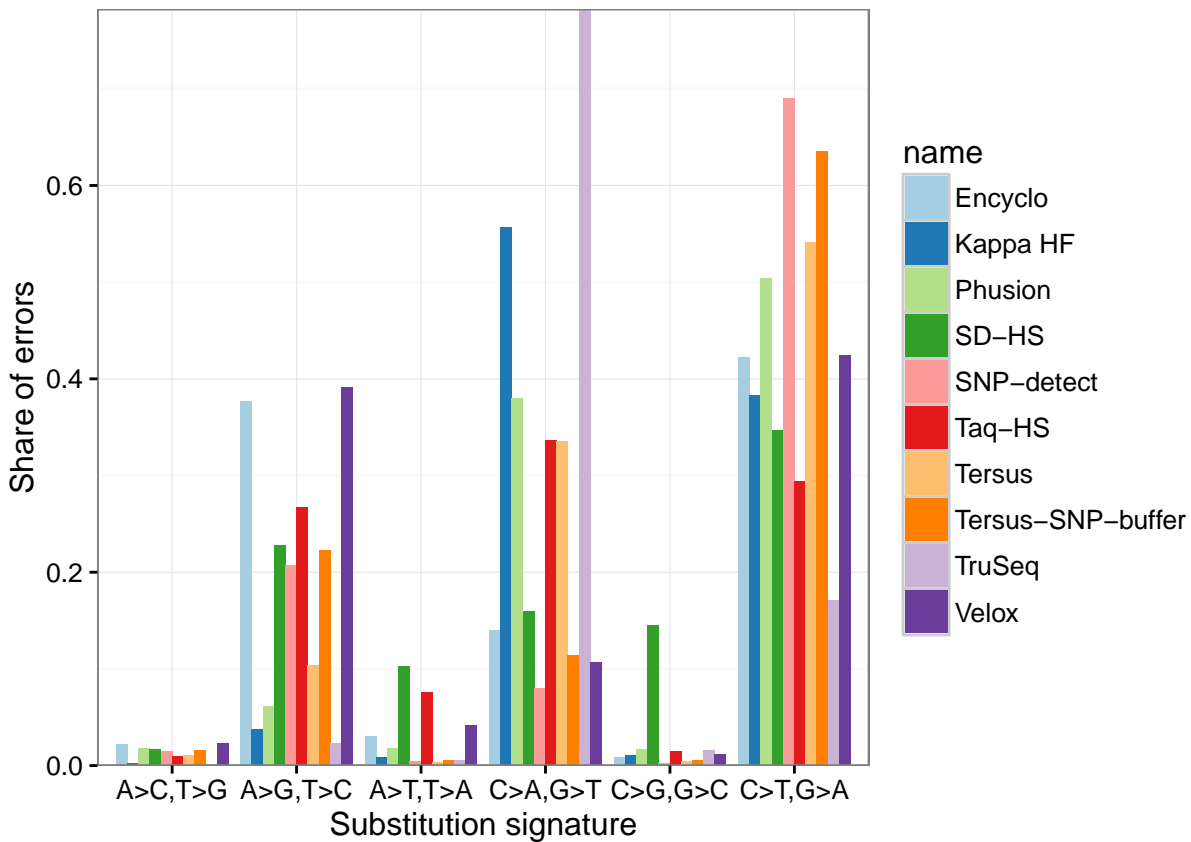
```

mutation.signature.rep = c("A>C,T>G","A>G,T>C","A>T,T>A","C>A,G>T","C>G,G>C","C>
df <- merge(df, sign.rep, all.x=T, all.y=F)

df.pattern <- ddply(df, .(name, mutation.signature.rep), summarize, count.sum = sum(count.major))
df.pattern <- ddply(df.pattern, .(name), transform, freq = count.sum / sum(count.sum))

ggplot(df.pattern, aes(x = mutation.signature.rep, weight = freq,
  fill = name)) + geom_bar(position = position_dodge()) +
  xlab("Substitution signature") + scale_y_continuous("Share of errors", expand=c(0,0)) +
  scale_fill_brewer(palette = "Paired") + theme_bw()

```



```

df.pattern.mat <- dcast(df.pattern, name ~ mutation.signature.rep, value.var = "freq")
rownames(df.pattern.mat) <- df.pattern.mat$name
df.pattern.mat$name <- NULL

df.pattern.mat <- as.matrix(df.pattern.mat)

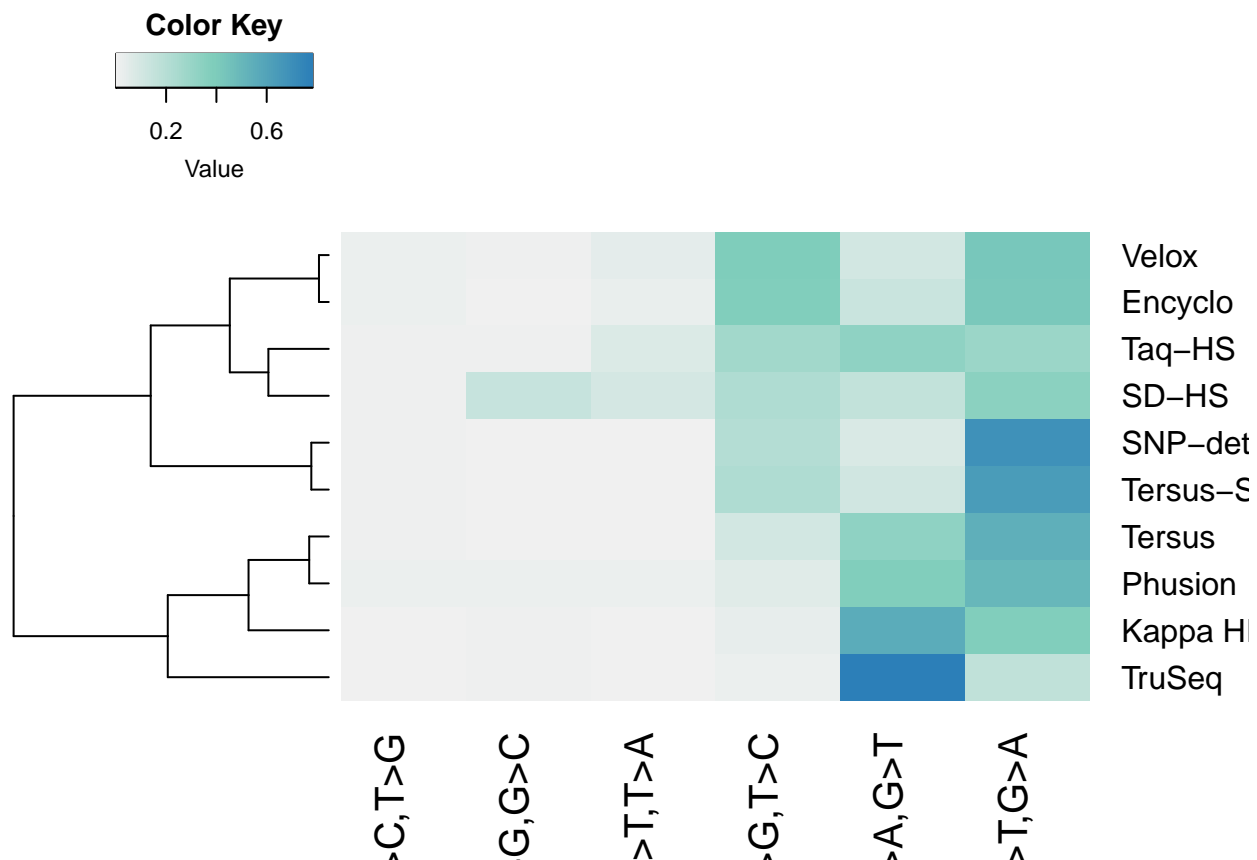
heatmap.2(df.pattern.mat, col=colorpanel(100, "#f0f0f0", "#7fcdbb", "#2c7fb8"),
  hclustfun = function(d) hclust(d, method="ward"),
  dendrogram = "row",
  density.info = "none", trace="none")

```

```

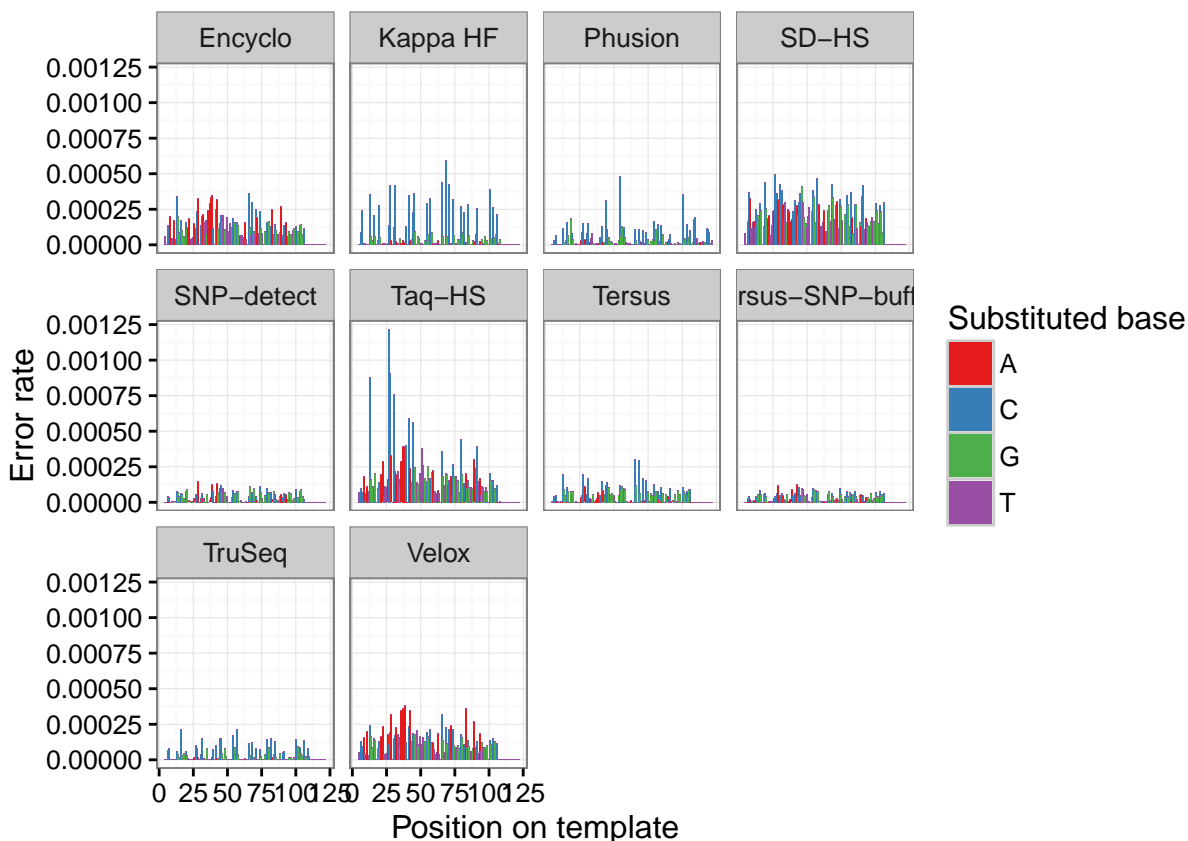
## The "ward" method has been renamed to "ward.D"; note new "ward.D2"
## The "ward" method has been renamed to "ward.D"; note new "ward.D2"

```



This goes to supplementary. Clear preference of substitution patterns for different polymerases, but no position-related trend.

```
ggplot(df, aes(x = mutation.pos, weight = count.major / coverage, fill = mutation.from)) +
  geom_histogram(bins = nchar(template)) + scale_fill_brewer("Substituted base", palette = "Set1") +
  xlab("Position on template") + ylab("Error rate") +
  facet_wrap(~name) + theme_bw()#, scales = "free_y") + theme_bw()
```



```
a <- aov(count.major / coverage ~ mutation.from * name + mutation.pos, df)
summary(a)
```

```
##               Df    Sum Sq   Mean Sq F value    Pr(>F)
## mutation.from    3 5.840e-07 1.946e-07  44.781 < 2e-16 ***
## name             9 1.064e-06 1.182e-07  27.207 < 2e-16 ***
## mutation.pos      1 1.000e-08 1.036e-08   2.384  0.123
## mutation.from:name 27 5.010e-07 1.857e-08   4.273 1.72e-12 ***
## Residuals       2121 9.217e-06 4.350e-09
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

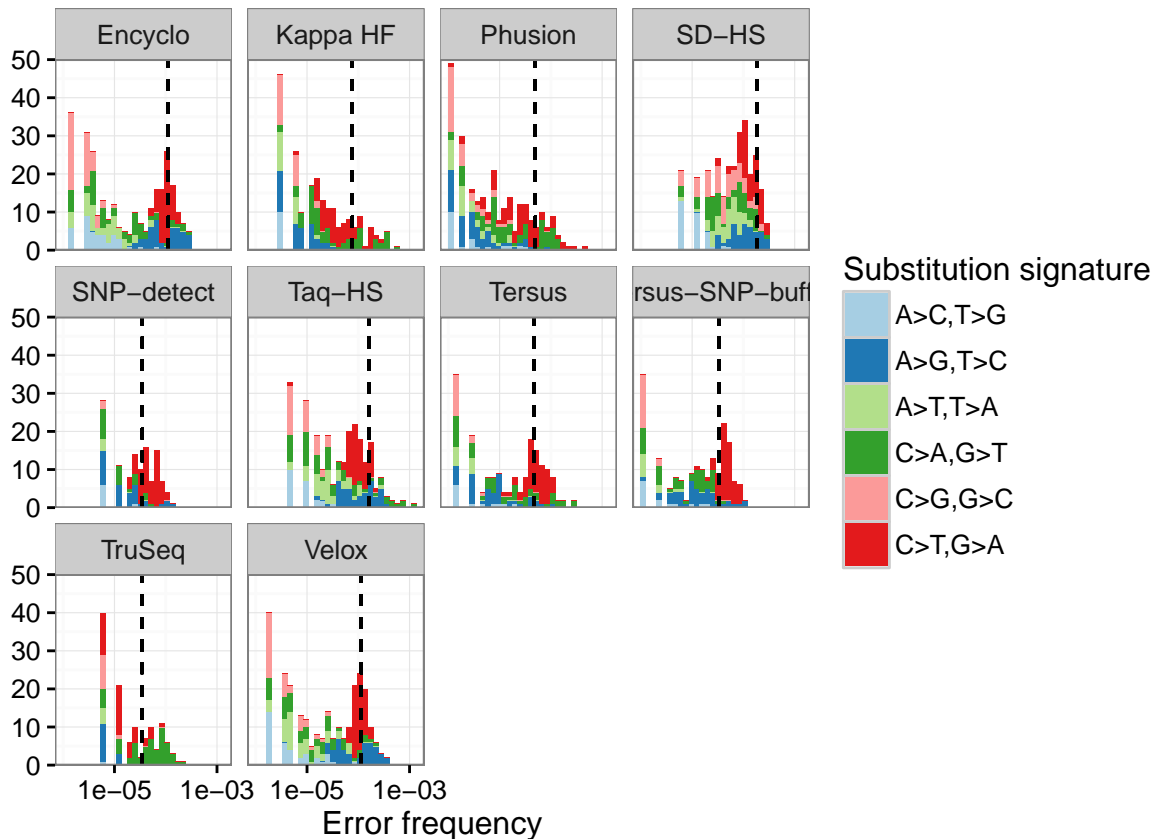
Error hotspot context pattern

```
df.mean.err <- ddply(df, .(name), summarize, mean.err.rate = sum(count.major) / mean(coverage) / nchar(
df <- merge(df, df.mean.err, all.x=T, all.y=F)

ggplot(df) +
  geom_histogram(aes(x = count.major / coverage, fill = mutation.signature.rep)) +
  geom_linerange(aes(x = mean.err.rate, ymin = 0, ymax=50), linetype = "dashed", color="black") +
  scale_fill_brewer("Substitution signature", palette = "Paired") +
  scale_x_log10("Error frequency") +
  scale_y_continuous("", expand=c(0,0)) + facet_wrap(~name) + theme_bw()
```



```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



```
#facet_wrap(~name, scales = "free_y") + theme_bw()
```

GC context profile

```
window.size <- 15
padding <- paste(rep("N", window.size), sep = "", collapse="")
template.p <- paste(padding, template, padding, sep="")

df$context <- sapply(df$mutation.pos,
  function(pos) substring(template.p, pos + 1, pos + 2 * window.size + 1))

df.context <- ddply(df, .(name, context), summarize,
  weight = sum(count.major))

context.rand <- sapply(sample(1:nchar(template), 1000, replace=T),
  function(pos) substring(template.p, pos + 1, pos + 2 * window.size + 1))

df.context.rand <- data.frame(context = context.rand)
df.context.rand$name <- "random"
df.context.rand$weight <- 1

df.context <- rbind(df.context, df.context.rand)
```

```

write.table(df.context, file = "context.txt", quote = F, sep = "\t", row.names = F)

system("groovy ProcessContext.groovy")

df.context.profile <- read.table("context.proc.txt", header=T, sep="\t")

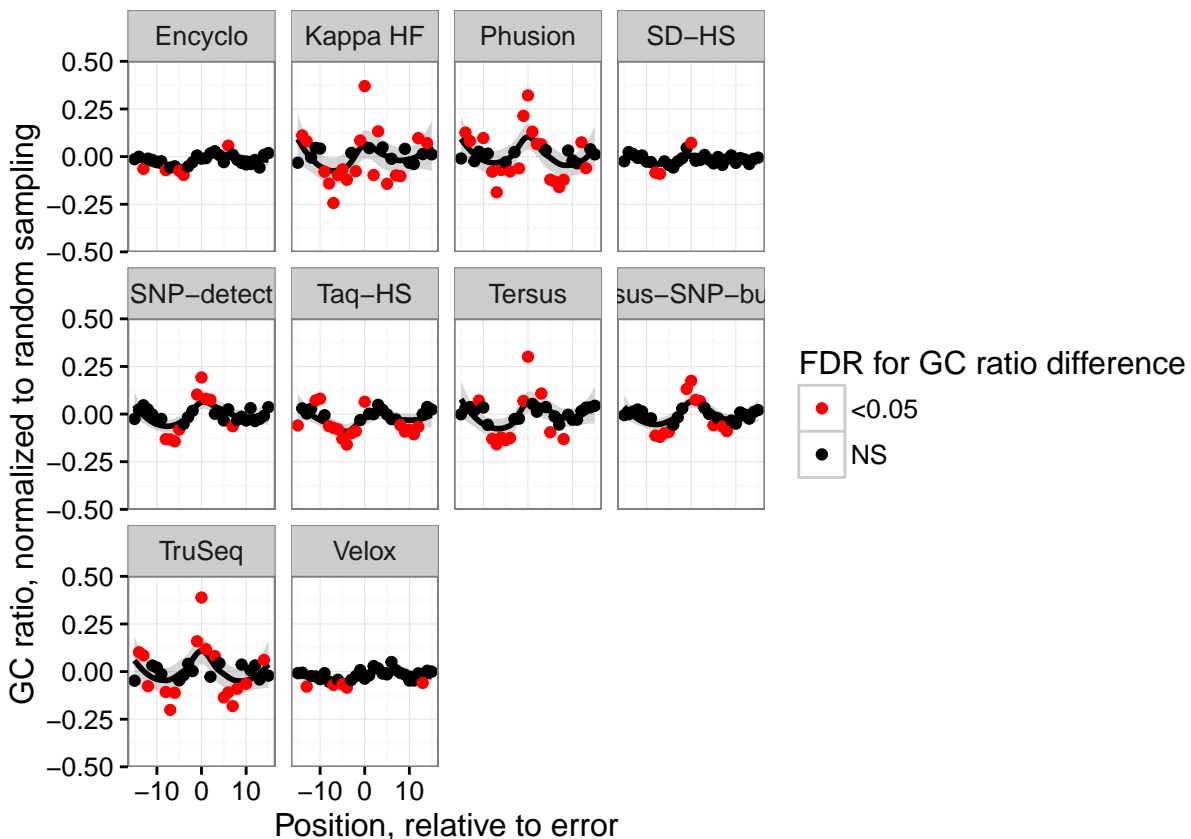
df.context.normalized <- merge(subset(df.context.profile, name != "random"),
                              subset(df.context.profile, name == "random"),
                              by = "pos")

df.context.normalized <- ddply(df.context.normalized, .(pos, name.x), transform,
                              pval = prop.test(x = value.y, n = sum.y, p = value.x / sum.x,
                                                alternative = "two.sided", correct = F)[[3]])

df.context.normalized$pval <- p.adjust(df.context.normalized$pval)

ggplot(df.context.normalized, aes(x = pos - window.size,
                                y = value.x / sum.x - value.y / sum.y)) +
  geom_smooth(colour="black") + geom_point(aes(color = factor(ifelse(pval < 0.05, "<0.05", "NS")))) +
  facet_wrap(~name.x) + scale_colour_manual("FDR for GC ratio difference", values=c("red", "black")) +
  scale_y_continuous("GC ratio, normalized to random sampling", limits = c(-0.5, 0.5), expand=c(0,0)) +
  theme_bw()

```



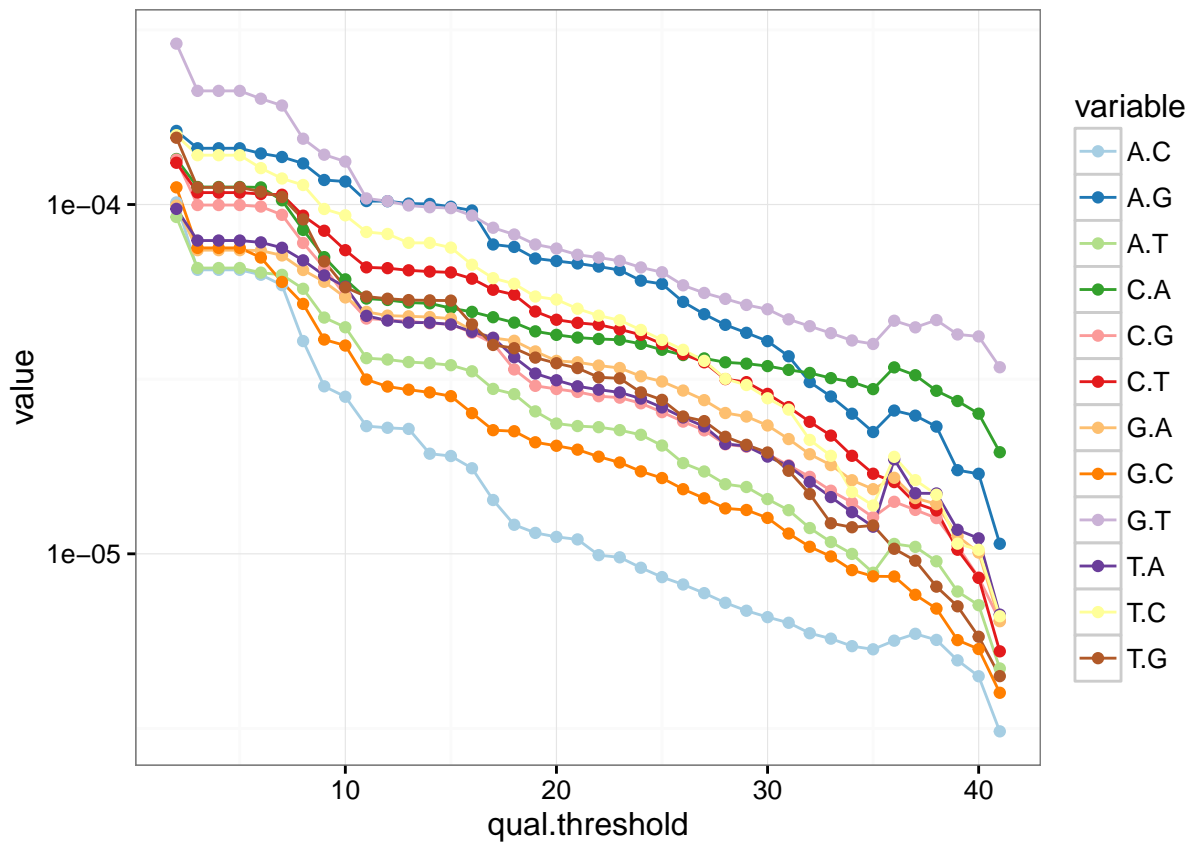
```

df.q <- read.table("mmqc.txt", header=T, sep="\t")
library(reshape2)

```

```
df.q <- melt(df.q, id.vars = "qual.threshold")

ggplot(df.q, aes(x=qual.threshold, color=variable, y=value)) +
  geom_line() + geom_point() + scale_y_log10() +
  scale_color_brewer(palette = "Paired") + theme_bw()
```



```
df.q$subset <- ifelse(df.q$qual.threshold >= 35, "high-quality", "low-quality")
df.q$signature <- ifelse(df.q$variable %in% c("G.T", "C.A"), "TrueSeq", "other")

#ggplot(df.q, aes(x=subset, group=interaction(subset, signature), fill=signature, y = value)) +
# geom_boxplot()
```