

BSH genes annotation (contigs)

Chunyu Zhao

11 December, 2018

Contents

1	20181208	1
1.1	taxonomic annotation	1
1.2	common ancestor	1
1.3	add bsh meta	1
1.4	per base coverage	1
2	20181209	1
2.1	bsh gene databases	1
2.2	bsh contigs set up pident and scov cutoff	1
2.2.1	(1) BSH distributuion for all raw blastx results	1
2.2.2	(2) filter threshold	4
2.2.3	(3) Keep top hit	4
2.3	common ancestor	4
2.4	metadata	4
3	20181210	5
3.1	bysample	5
3.1.1	all five time points linear mixed effects	6
3.1.2	formula effect at 4.5 month	7
3.2	byphylum	7

1 20181208

1.1 taxonomic annotation

1.2 common ancestor

1.3 add bsh meta

1.4 per base coverage

- 20181120: added `sum` from the `per.base.df` data frame to `sunbeam/coverage.rules`.

2 20181209

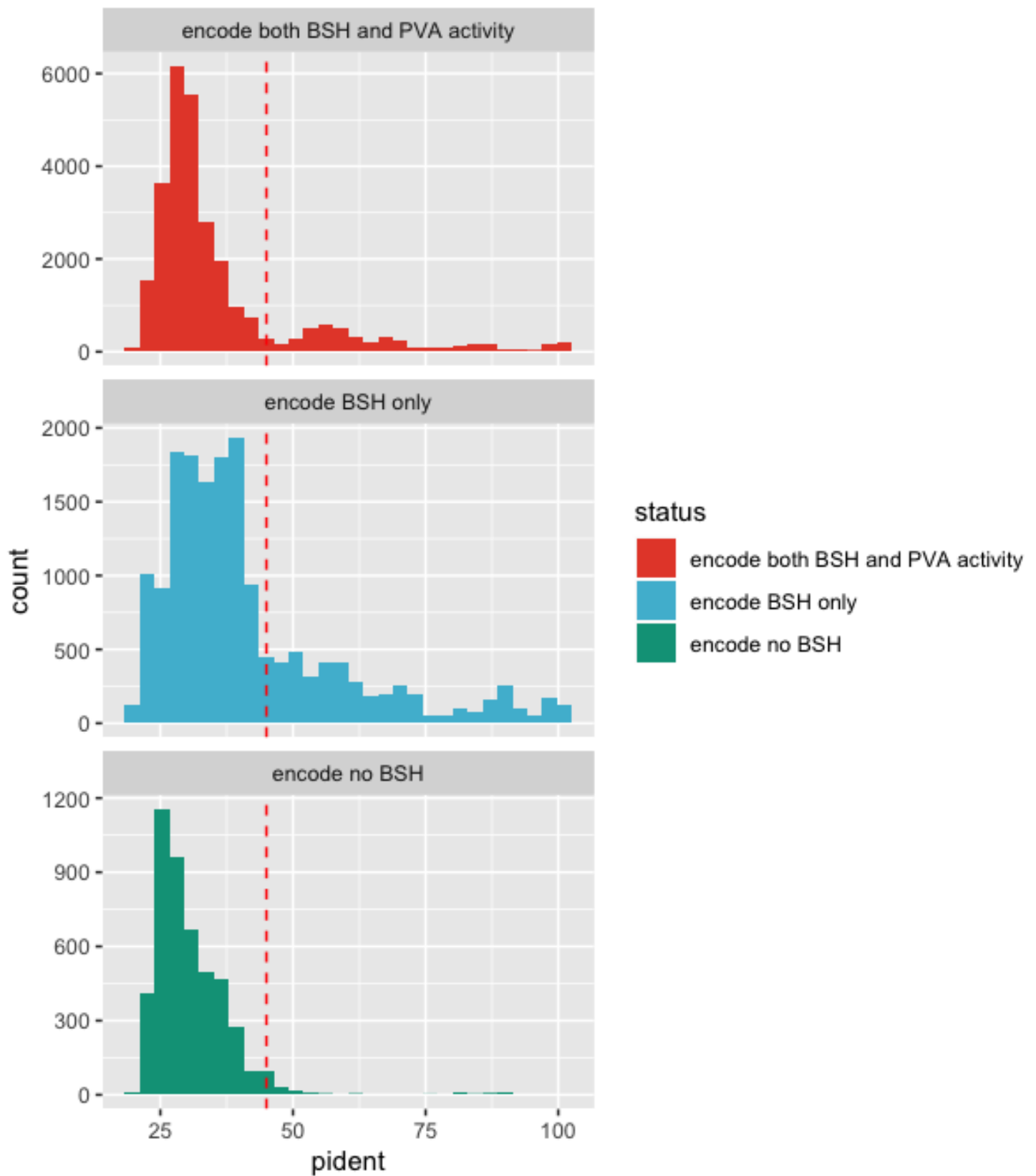
2.1 bsh gene databases

2.2 bsh contigs set up pident and scov cutoff

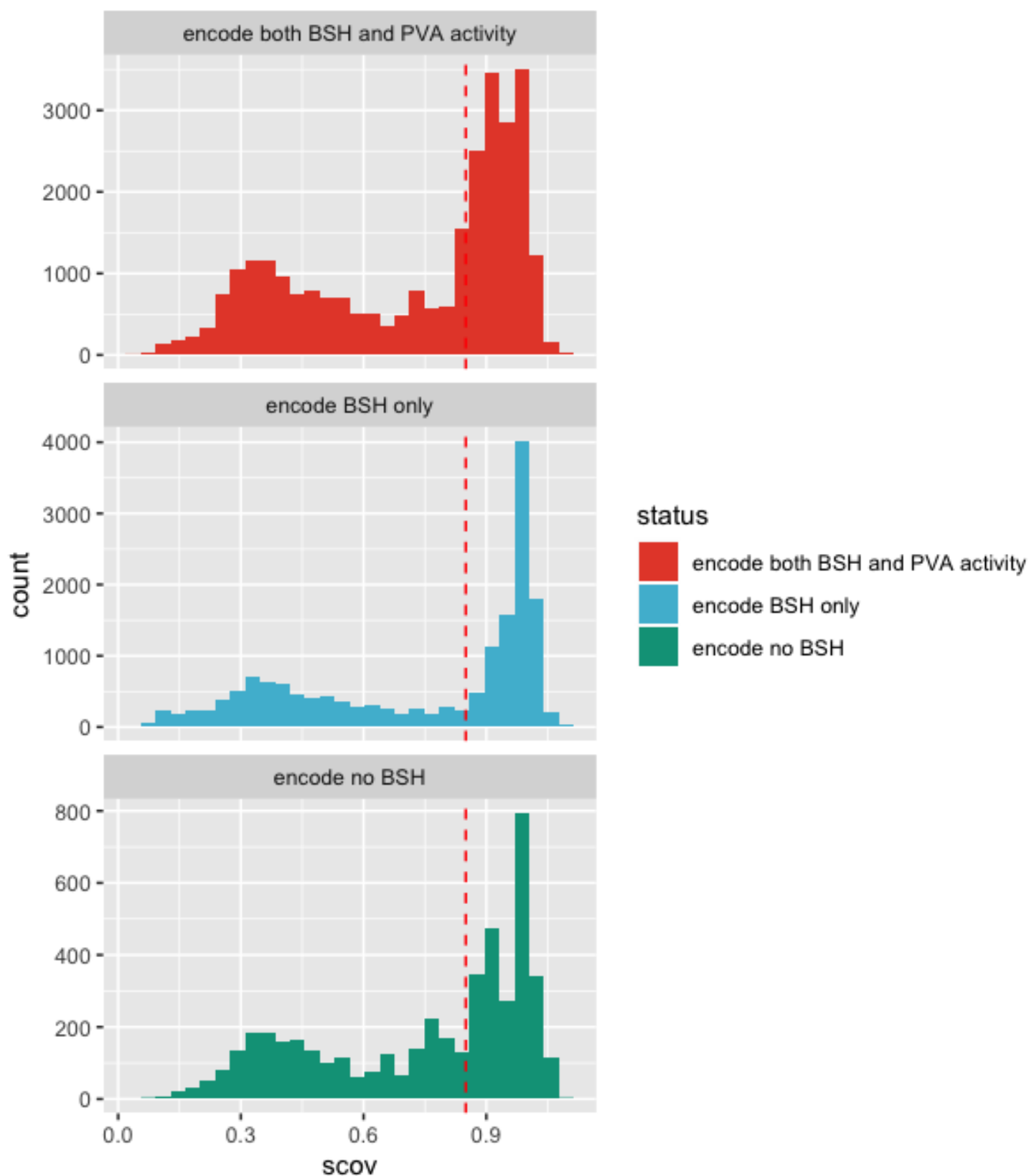
2.2.1 (1) BSH distributuion for all raw blastx results

- **task:** detect whether out shotgun data encode BSH genes or not; don't care species encode that genes for now.

distribution for ident



distribution for bsh gene coverage



2.2.2 (2) filter threshold

- set `pident.cutoff` to 45 and `scov.cutoff` to 0.85 for the downstream analysis.

2.2.3 (3) Keep top hit

Since we are only interested in the presence/absence of the BSH genes, it works towards our benefits if the one contig blasted to multiple BSH genes, meaning highly homologous genes. ****Unless***, the same contig mapped to bash **encode BSH** and **no encode BSH** groups.

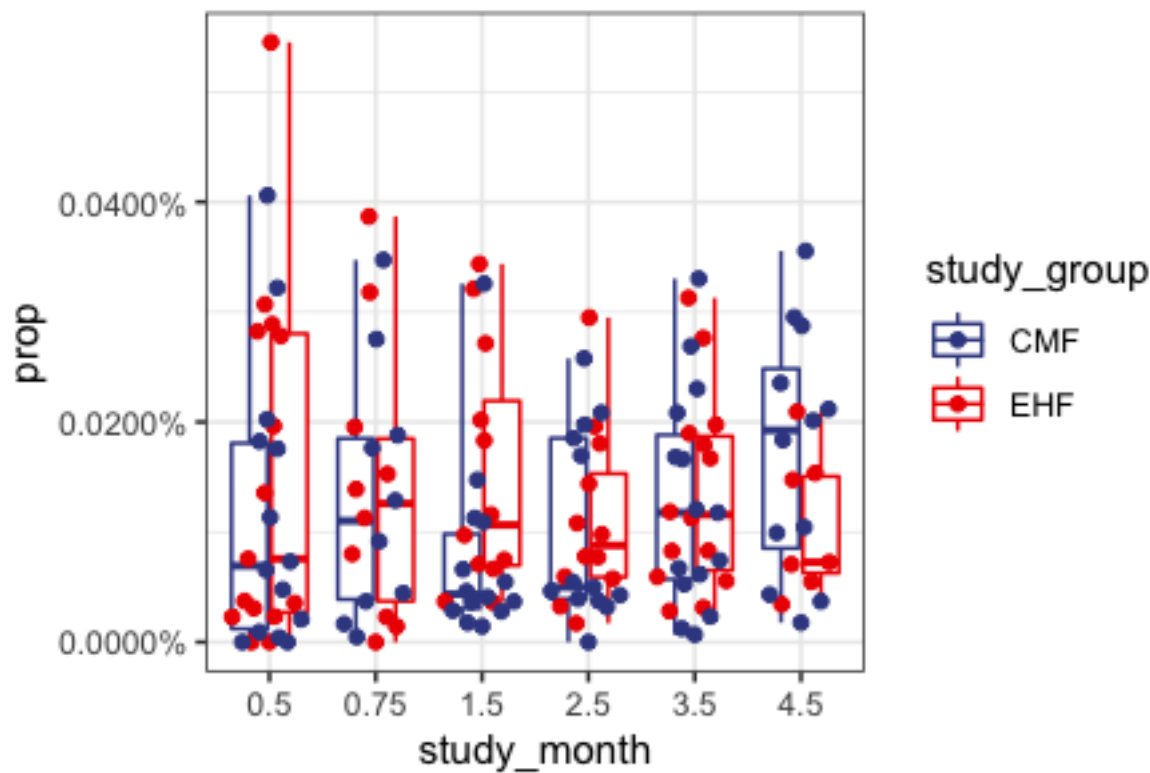
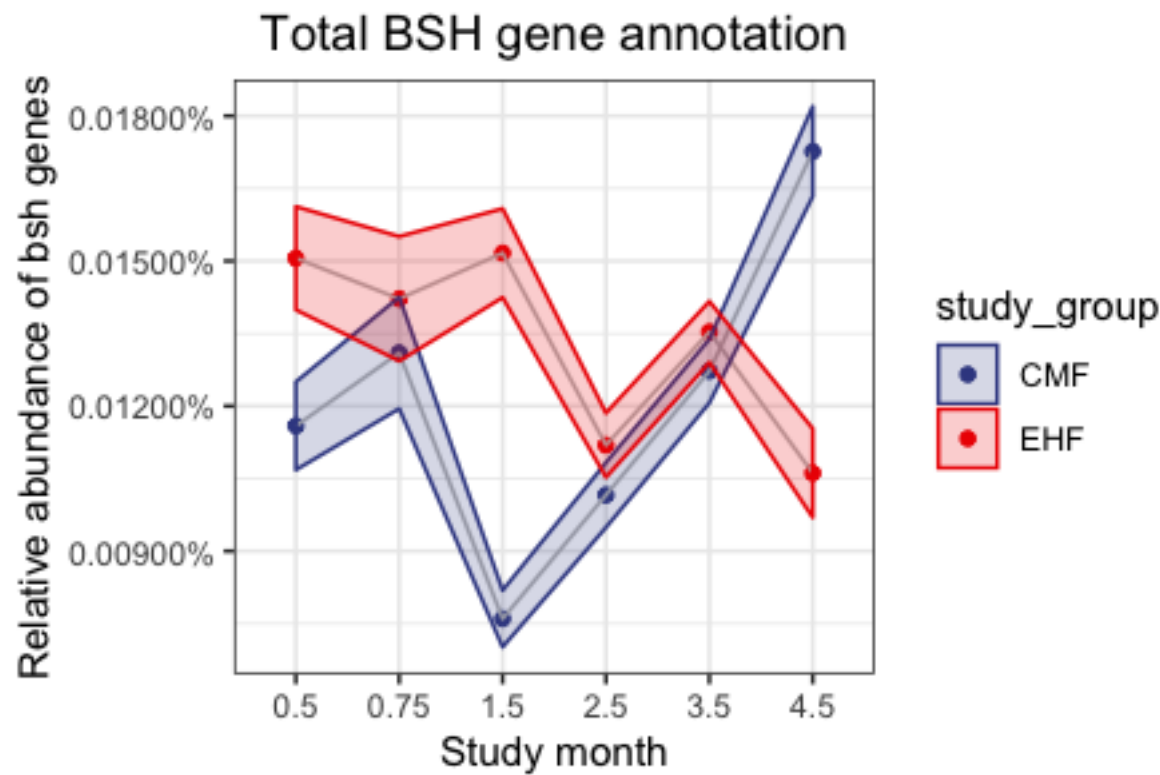
Table 1: contigs with paralogs, need to be careful

sample	qseqid	status
GRO_009	k141_23138	encode both BSH and PVA activity
GRO_009	k141_23138	encode no BSH
GRO_010	k141_54210	encode both BSH and PVA activity
GRO_010	k141_54210	encode no BSH
GRO_087	k141_14049	encode both BSH and PVA activity
GRO_087	k141_14049	encode no BSH
GRO_097	k141_9432	encode no BSH
GRO_097	k141_9432	encode both BSH and PVA activity
GRO_164	k141_14599	encode both BSH and PVA activity
GRO_164	k141_14599	encode no BSH

2.3 common ancestor

2.4 metadata

20181112: Also, just a reminder - both 0.5 and 0.75 mo timepoints are baseline.



3.1.1 all five time points linear mixed effects

```
## Linear mixed-effects model fit by REML
## Data: bysample
##      AIC      BIC    logLik
## 272.8918 313.669 -122.4459
##
## Random effects:
## Formula: ~1 | SubjectID
##      (Intercept) Residual
## StdDev:  0.3106412 0.475512
##
## Fixed effects: LogProp ~ study_group * study_month
##
##              Value Std.Error DF   t-value p-value
## (Intercept)   -4.379732 0.1508499 108 -29.033708 0.0000
## study_groupEHF    0.133763 0.2103878  28  0.635793 0.5301
## study_month0.75    0.283551 0.2005191 108  1.414084 0.1602
## study_month1.5     0.125249 0.1797267 108  0.696888 0.4874
## study_month2.5     0.116977 0.1846394 108  0.633541 0.5277
## study_month3.5     0.305756 0.1771909 108  1.725576 0.0873
## study_month4.5     0.486450 0.1892276 108  2.570716 0.0115
## study_groupEHF:study_month0.75 -0.204189 0.2812770 108 -0.725934 0.4694
## study_groupEHF:study_month1.5  0.206775 0.2586372 108  0.799479 0.4258
## study_groupEHF:study_month2.5  0.056324 0.2618975 108  0.215061 0.8301
## study_groupEHF:study_month3.5 -0.006894 0.2505910 108 -0.027510 0.9781
## study_groupEHF:study_month4.5 -0.251939 0.2931081 108 -0.859542 0.3919
## Correlation:
##              (Intr) st_EHF s_0.75 st_1.5 st_2.5 st_3.5
## study_groupEHF   -0.717
## study_month0.75  -0.543  0.389
## study_month1.5   -0.596  0.427  0.448
## study_month2.5   -0.587  0.421  0.439  0.487
## study_month3.5   -0.611  0.438  0.462  0.507  0.500
## study_month4.5   -0.573  0.411  0.434  0.475  0.467  0.488
## study_groupEHF:study_month0.75  0.387 -0.532 -0.713 -0.319 -0.313 -0.329
## study_groupEHF:study_month1.5  0.414 -0.574 -0.311 -0.695 -0.338 -0.352
## study_groupEHF:study_month2.5  0.414 -0.570 -0.310 -0.343 -0.705 -0.352
## study_groupEHF:study_month3.5  0.432 -0.596 -0.327 -0.359 -0.353 -0.707
## study_groupEHF:study_month4.5  0.370 -0.510 -0.280 -0.307 -0.301 -0.315
##              st_4.5 s_EHF:_0 s_EHF:_1 s_EHF:_2 s_EHF:_3
## study_groupEHF
## study_month0.75
## study_month1.5
## study_month2.5
## study_month3.5
## study_month4.5
## study_groupEHF:study_month0.75 -0.310
## study_groupEHF:study_month1.5 -0.330  0.434
## study_groupEHF:study_month2.5 -0.329  0.421  0.458
## study_groupEHF:study_month3.5 -0.345  0.445  0.481  0.482
## study_groupEHF:study_month4.5 -0.646  0.386  0.421  0.399  0.425
##
## Standardized Within-Group Residuals:
##      Min      Q1      Med      Q3      Max
## -3.085656 -0.357015 0.152381 0.482992 1.842638
##
```

```
## Number of Observations: 148
## Number of Groups: 30
```

3.1.2 formula effect at 4.5 month

```
##
## Call:
## lm(formula = LogProp ~ study_group, data = s_toTest)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.8364 -0.1580  0.1599  0.2495  0.4455
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -3.8935     0.1065  -36.57  <2e-16 ***
## study_groupEHF -0.1490     0.1754   -0.85   0.407
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.3688 on 17 degrees of freedom
## Multiple R-squared:  0.04074,    Adjusted R-squared:  -0.01569
## F-statistic: 0.7219 on 1 and 17 DF,  p-value: 0.4073
```

3.2 byphylum



Total BSH gene annotation

