

Supplementary material for UniqTag: Content-derived unique and stable identifiers for gene annotation

Shaun Jackman

Supplementary material

The following supplementary material of the UniqTag paper give the R code and the data, shown in supplementary table S1, used to generate figure 1 from the main material of the paper and supplementary figure S1.

Load libraries

```
library(ggplot2)
library(knitr) # for kable
library(reshape2)
library(scales) # for alpha
```

Read the data

```
data.orig <- read.delim('uniqtag.tsv',
  colClasses = c(A = 'factor', B = 'factor'))
x <- do.call(rbind, strsplit(as.character(data.orig$Table), '.', fixed = TRUE))
colnames(x) <- c('Data', 'Transform', 'Identifier')
data <- cbind(data.orig, x)
rm(x)
data$k <- as.integer(gsub('[a-z]*', '', data$Identifier))

build.wide <- with(data,
  data.frame(Build.A = A, Build.B = B,
    Num.A = Only.A + Both, Num.B = Only.B + Both))
build.tall <- melt(build.wide, id.vars = c('Build.A', 'Build.B'),
  variable.name = 'Build', value.name = 'Count')
```

Figure 1. Plot the number of common identifiers vs. older build

The number of common UniqTag identifiers between older builds of the Ensembl human genome and the current build 75, the number of common gene and protein identifiers between builds, and the number of genes with peptide sequences that are identical between builds.

```
data.subset <- subset(data, data$k == 9 | is.na(data$k))
aes.data <- aes(x = A, y = Both,
  group = Table, colour = Identifier)
aes.build <- aes(x = Build.A, y = Count,
```

```

    group = Build, linetype = Build, shape = Build)
ggplot() +
  geom_point(aes.data, data.subset) +
  geom_line(aes.data, data.subset) +
  scale_colour_brewer(palette = 'Set1',
    breaks = c('gene', 'uniquetag9', 'id', 'seq'),
    labels = c('Gene ID (ENSG)', 'UniqTag (9-mer)',
      'Protein ID (ENSP)', 'Identical peptide sequence')) +

  geom_point(aes.build, build.tall) +
  geom_line(aes.build, build.tall) +
  scale_linetype_manual(name = 'Number of genes',
    breaks = c('Num.B', 'Num.A'),
    labels = c('Ensembl build 75', 'Older Ensembl build'),
    values = c('solid', 'dashed')) +
  scale_shape_manual(name = 'Number of genes',
    breaks = c('Num.B', 'Num.A'),
    labels = c('Ensembl build 75', 'Older Ensembl build'),
    values = c(20, 32)) +

  theme_bw() +
  theme(legend.position = c(1.0, 0),
    legend.justification = c(1, 0),
    legend.box.just = 'right',
    legend.background = element_rect(fill = alpha('white', 0))) +
  xlab('Older Ensembl build') +
  ylab('Identifiers in common with Ensembl build 75')

```

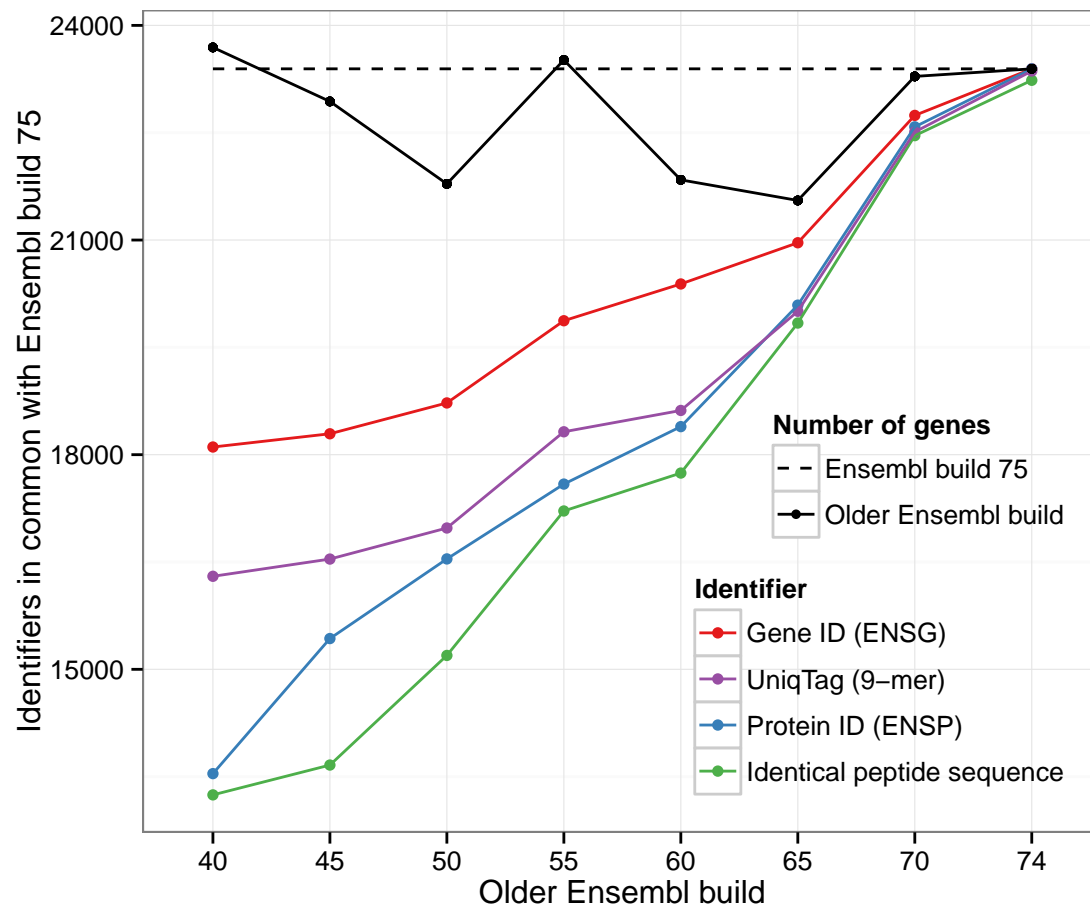


Figure S1. Plot the number of common identifiers vs. k

The number of common UniqTag identifiers between older builds of the Ensembl human genome and the current build 75 for different values of k.

```
ggplot(na.omit(data), aes(x = k, y = Both, group = A, colour = A)) +
  geom_point() +
  geom_line() +
  scale_x_continuous(trans = log_trans(),
    breaks = c(1, 2, 5, 10, 20, 50, 100, 200)) +
  scale_colour_brewer(palette = 'Set2') +
  guides(colour = guide_legend(reverse = TRUE)) +
  theme_bw() +
  xlab('Size of UniqTag k-mer (bp)') +
  ylab('Identifiers in common with Ensembl build 75')
```

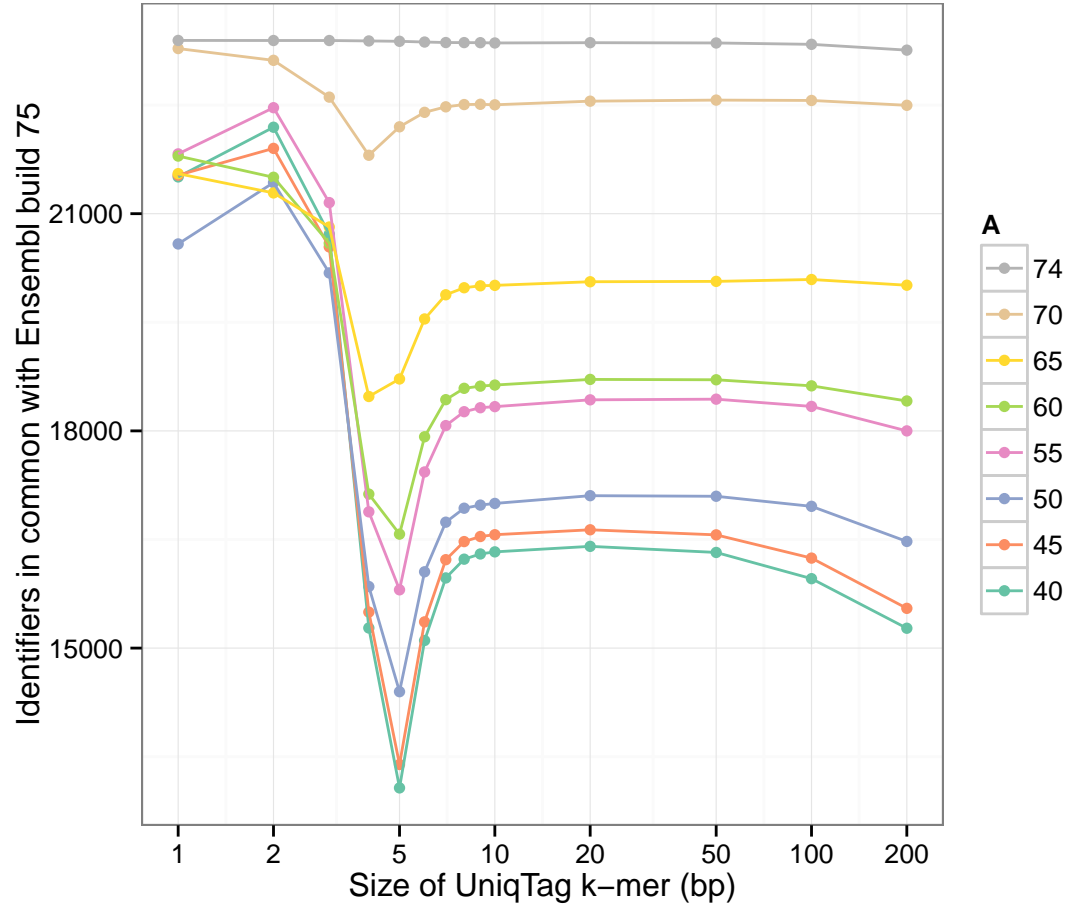


Table S1. The number of common identifiers

The number of common identifiers between older builds of the Ensembl human genome and the current build 75 for gene ID (ENSG), protein ID (ENSP), exact peptide sequence and UniqTag for different values of k . These data are used to plot the above figures. It is also available in tab-separated values (TSV) format.

`kable(data)`

Table	A	B	Only.A	Both	Only.B	Data	Transform	Identifier	k
all.uniqlgene.gene	40	75	5585	18107	5286	all	uniqlgene	gene	NA
all.uniqlgene.gene	45	75	4645	18292	5101	all	uniqlgene	gene	NA
all.uniqlgene.gene	50	75	3062	18723	4670	all	uniqlgene	gene	NA
all.uniqlgene.gene	55	75	3644	19872	3521	all	uniqlgene	gene	NA
all.uniqlgene.gene	60	75	1455	20386	3007	all	uniqlgene	gene	NA
all.uniqlgene.gene	65	75	591	20962	2431	all	uniqlgene	gene	NA
all.uniqlgene.gene	70	75	545	22742	651	all	uniqlgene	gene	NA
all.uniqlgene.gene	74	75	0	23393	0	all	uniqlgene	gene	NA
all.uniqlgene.id	40	75	10150	13542	9851	all	uniqlgene	id	NA

Table	A	B	Only.A	Both	Only.B	Data	Transform	Identifier	k
all.uniqgene.id	45	75	7507	15430	7963	all	uniqgene	id	NA
all.uniqgene.id	50	75	5242	16543	6850	all	uniqgene	id	NA
all.uniqgene.id	55	75	5927	17589	5804	all	uniqgene	id	NA
all.uniqgene.id	60	75	3449	18392	5001	all	uniqgene	id	NA
all.uniqgene.id	65	75	1463	20090	3303	all	uniqgene	id	NA
all.uniqgene.id	70	75	705	22582	811	all	uniqgene	id	NA
all.uniqgene.id	74	75	0	23393	0	all	uniqgene	id	NA
all.uniqgene.seq	40	75	10447	13245	10148	all	uniqgene	seq	NA
all.uniqgene.seq	45	75	9275	13662	9731	all	uniqgene	seq	NA
all.uniqgene.seq	50	75	6591	15194	8199	all	uniqgene	seq	NA
all.uniqgene.seq	55	75	6303	17213	6180	all	uniqgene	seq	NA
all.uniqgene.seq	60	75	4098	17743	5650	all	uniqgene	seq	NA
all.uniqgene.seq	65	75	1713	19840	3553	all	uniqgene	seq	NA
all.uniqgene.seq	70	75	828	22459	934	all	uniqgene	seq	NA
all.uniqgene.seq	74	75	160	23233	160	all	uniqgene	seq	NA
all.uniqgene.uniqttag1	40	75	2184	21508	1885	all	uniqgene	uniqttag1	1
all.uniqgene.uniqttag1	45	75	1405	21532	1861	all	uniqgene	uniqttag1	1
all.uniqgene.uniqttag1	50	75	1203	20582	2811	all	uniqgene	uniqttag1	1
all.uniqgene.uniqttag1	55	75	1690	21826	1567	all	uniqgene	uniqttag1	1
all.uniqgene.uniqttag1	60	75	45	21796	1597	all	uniqgene	uniqttag1	1
all.uniqgene.uniqttag1	65	75	0	21553	1840	all	uniqgene	uniqttag1	1
all.uniqgene.uniqttag1	70	75	6	23281	112	all	uniqgene	uniqttag1	1
all.uniqgene.uniqttag1	74	75	0	23393	0	all	uniqgene	uniqttag1	1
all.uniqgene.uniqttag2	40	75	1498	22194	1199	all	uniqgene	uniqttag2	2
all.uniqgene.uniqttag2	45	75	1035	21902	1491	all	uniqgene	uniqttag2	2
all.uniqgene.uniqttag2	50	75	356	21429	1964	all	uniqgene	uniqttag2	2
all.uniqgene.uniqttag2	55	75	1052	22464	929	all	uniqgene	uniqttag2	2
all.uniqgene.uniqttag2	60	75	338	21503	1890	all	uniqgene	uniqttag2	2
all.uniqgene.uniqttag2	65	75	266	21287	2106	all	uniqgene	uniqttag2	2
all.uniqgene.uniqttag2	70	75	169	23118	275	all	uniqgene	uniqttag2	2
all.uniqgene.uniqttag2	74	75	1	23392	1	all	uniqgene	uniqttag2	2
all.uniqgene.uniqttag3	40	75	2975	20717	2676	all	uniqgene	uniqttag3	3
all.uniqgene.uniqttag3	45	75	2396	20541	2852	all	uniqgene	uniqttag3	3
all.uniqgene.uniqttag3	50	75	1603	20182	3211	all	uniqgene	uniqttag3	3
all.uniqgene.uniqttag3	55	75	2363	21153	2240	all	uniqgene	uniqttag3	3
all.uniqgene.uniqttag3	60	75	1249	20592	2801	all	uniqgene	uniqttag3	3

Table	A	B	Only.A	Both	Only.B	Data	Transform	Identifier	k
all.uniqlgene.uniqltag3	65	75	737	20816	2577	all	uniqlgene	uniqltag3	3
all.uniqlgene.uniqltag3	70	75	677	22610	783	all	uniqlgene	uniqltag3	3
all.uniqlgene.uniqltag3	74	75	1	23392	1	all	uniqlgene	uniqltag3	3
all.uniqlgene.uniqltag4	40	75	8414	15278	8115	all	uniqlgene	uniqltag4	4
all.uniqlgene.uniqltag4	45	75	7440	15497	7896	all	uniqlgene	uniqltag4	4
all.uniqlgene.uniqltag4	50	75	5935	15850	7543	all	uniqlgene	uniqltag4	4
all.uniqlgene.uniqltag4	55	75	6634	16882	6511	all	uniqlgene	uniqltag4	4
all.uniqlgene.uniqltag4	60	75	4714	17127	6266	all	uniqlgene	uniqltag4	4
all.uniqlgene.uniqltag4	65	75	3078	18475	4918	all	uniqlgene	uniqltag4	4
all.uniqlgene.uniqltag4	70	75	1480	21807	1586	all	uniqlgene	uniqltag4	4
all.uniqlgene.uniqltag4	74	75	7	23386	7	all	uniqlgene	uniqltag4	4
all.uniqlgene.uniqltag5	40	75	10623	13069	10324	all	uniqlgene	uniqltag5	5
all.uniqlgene.uniqltag5	45	75	9545	13392	10001	all	uniqlgene	uniqltag5	5
all.uniqlgene.uniqltag5	50	75	7387	14398	8995	all	uniqlgene	uniqltag5	5
all.uniqlgene.uniqltag5	55	75	7711	15805	7588	all	uniqlgene	uniqltag5	5
all.uniqlgene.uniqltag5	60	75	5267	16574	6819	all	uniqlgene	uniqltag5	5
all.uniqlgene.uniqltag5	65	75	2836	18717	4676	all	uniqlgene	uniqltag5	5
all.uniqlgene.uniqltag5	70	75	1087	22200	1193	all	uniqlgene	uniqltag5	5
all.uniqlgene.uniqltag5	74	75	12	23381	12	all	uniqlgene	uniqltag5	5
all.uniqlgene.uniqltag6	40	75	8587	15105	8288	all	uniqlgene	uniqltag6	6
all.uniqlgene.uniqltag6	45	75	7575	15362	8031	all	uniqlgene	uniqltag6	6
all.uniqlgene.uniqltag6	50	75	5731	16054	7339	all	uniqlgene	uniqltag6	6
all.uniqlgene.uniqltag6	55	75	6083	17433	5960	all	uniqlgene	uniqltag6	6
all.uniqlgene.uniqltag6	60	75	3922	17919	5474	all	uniqlgene	uniqltag6	6
all.uniqlgene.uniqltag6	65	75	2007	19546	3847	all	uniqlgene	uniqltag6	6
all.uniqlgene.uniqltag6	70	75	887	22400	993	all	uniqlgene	uniqltag6	6
all.uniqlgene.uniqltag6	74	75	22	23371	22	all	uniqlgene	uniqltag6	6
all.uniqlgene.uniqltag7	40	75	7723	15969	7424	all	uniqlgene	uniqltag7	7
all.uniqlgene.uniqltag7	45	75	6716	16221	7172	all	uniqlgene	uniqltag7	7
all.uniqlgene.uniqltag7	50	75	5046	16739	6654	all	uniqlgene	uniqltag7	7
all.uniqlgene.uniqltag7	55	75	5443	18073	5320	all	uniqlgene	uniqltag7	7
all.uniqlgene.uniqltag7	60	75	3410	18431	4962	all	uniqlgene	uniqltag7	7
all.uniqlgene.uniqltag7	65	75	1673	19880	3513	all	uniqlgene	uniqltag7	7
all.uniqlgene.uniqltag7	70	75	811	22476	917	all	uniqlgene	uniqltag7	7
all.uniqlgene.uniqltag7	74	75	29	23364	29	all	uniqlgene	uniqltag7	7
all.uniqlgene.uniqltag8	40	75	7464	16228	7165	all	uniqlgene	uniqltag8	8

Table	A	B	Only.A	Both	Only.B	Data	Transform	Identifier	k
all.uniqlgene.uniqltag8	45	75	6466	16471	6922	all	uniqlgene	uniqltag8	8
all.uniqlgene.uniqltag8	50	75	4853	16932	6461	all	uniqlgene	uniqltag8	8
all.uniqlgene.uniqltag8	55	75	5251	18265	5128	all	uniqlgene	uniqltag8	8
all.uniqlgene.uniqltag8	60	75	3253	18588	4805	all	uniqlgene	uniqltag8	8
all.uniqlgene.uniqltag8	65	75	1576	19977	3416	all	uniqlgene	uniqltag8	8
all.uniqlgene.uniqltag8	70	75	780	22507	886	all	uniqlgene	uniqltag8	8
all.uniqlgene.uniqltag8	74	75	30	23363	30	all	uniqlgene	uniqltag8	8
all.uniqlgene.uniqltag9	40	75	7392	16300	7093	all	uniqlgene	uniqltag9	9
all.uniqlgene.uniqltag9	45	75	6396	16541	6852	all	uniqlgene	uniqltag9	9
all.uniqlgene.uniqltag9	50	75	4810	16975	6418	all	uniqlgene	uniqltag9	9
all.uniqlgene.uniqltag9	55	75	5196	18320	5073	all	uniqlgene	uniqltag9	9
all.uniqlgene.uniqltag9	60	75	3223	18618	4775	all	uniqlgene	uniqltag9	9
all.uniqlgene.uniqltag9	65	75	1549	20004	3389	all	uniqlgene	uniqltag9	9
all.uniqlgene.uniqltag9	70	75	776	22511	882	all	uniqlgene	uniqltag9	9
all.uniqlgene.uniqltag9	74	75	31	23362	31	all	uniqlgene	uniqltag9	9
all.uniqlgene.uniqltag10	40	75	7363	16329	7064	all	uniqlgene	uniqltag10	10
all.uniqlgene.uniqltag10	45	75	6371	16566	6827	all	uniqlgene	uniqltag10	10
all.uniqlgene.uniqltag10	50	75	4787	16998	6395	all	uniqlgene	uniqltag10	10
all.uniqlgene.uniqltag10	55	75	5181	18335	5058	all	uniqlgene	uniqltag10	10
all.uniqlgene.uniqltag10	60	75	3208	18633	4760	all	uniqlgene	uniqltag10	10
all.uniqlgene.uniqltag10	65	75	1543	20010	3383	all	uniqlgene	uniqltag10	10
all.uniqlgene.uniqltag10	70	75	783	22504	889	all	uniqlgene	uniqltag10	10
all.uniqlgene.uniqltag10	74	75	35	23358	35	all	uniqlgene	uniqltag10	10
all.uniqlgene.uniqltag20	40	75	7287	16405	6988	all	uniqlgene	uniqltag20	20
all.uniqlgene.uniqltag20	45	75	6303	16634	6759	all	uniqlgene	uniqltag20	20
all.uniqlgene.uniqltag20	50	75	4680	17105	6288	all	uniqlgene	uniqltag20	20
all.uniqlgene.uniqltag20	55	75	5087	18429	4964	all	uniqlgene	uniqltag20	20
all.uniqlgene.uniqltag20	60	75	3130	18711	4682	all	uniqlgene	uniqltag20	20
all.uniqlgene.uniqltag20	65	75	1493	20060	3333	all	uniqlgene	uniqltag20	20
all.uniqlgene.uniqltag20	70	75	733	22554	839	all	uniqlgene	uniqltag20	20
all.uniqlgene.uniqltag20	74	75	31	23362	31	all	uniqlgene	uniqltag20	20
all.uniqlgene.uniqltag50	40	75	7371	16321	7072	all	uniqlgene	uniqltag50	50
all.uniqlgene.uniqltag50	45	75	6373	16564	6829	all	uniqlgene	uniqltag50	50
all.uniqlgene.uniqltag50	50	75	4688	17097	6296	all	uniqlgene	uniqltag50	50
all.uniqlgene.uniqltag50	55	75	5078	18438	4955	all	uniqlgene	uniqltag50	50
all.uniqlgene.uniqltag50	60	75	3135	18706	4687	all	uniqlgene	uniqltag50	50

Table	A	B	Only.A	Both	Only.B	Data	Transform	Identifier	k
all.uniqlgene.uniqltag50	65	75	1488	20065	3328	all	uniqlgene	uniqltag50	50
all.uniqlgene.uniqltag50	70	75	718	22569	824	all	uniqlgene	uniqltag50	50
all.uniqlgene.uniqltag50	74	75	35	23358	35	all	uniqlgene	uniqltag50	50
all.uniqlgene.uniqltag100	40	75	7733	15959	7434	all	uniqlgene	uniqltag100	100
all.uniqlgene.uniqltag100	45	75	6694	16243	7150	all	uniqlgene	uniqltag100	100
all.uniqlgene.uniqltag100	50	75	4827	16958	6435	all	uniqlgene	uniqltag100	100
all.uniqlgene.uniqltag100	55	75	5178	18338	5055	all	uniqlgene	uniqltag100	100
all.uniqlgene.uniqltag100	60	75	3219	18622	4771	all	uniqlgene	uniqltag100	100
all.uniqlgene.uniqltag100	65	75	1462	20091	3302	all	uniqlgene	uniqltag100	100
all.uniqlgene.uniqltag100	70	75	723	22564	829	all	uniqlgene	uniqltag100	100
all.uniqlgene.uniqltag100	74	75	54	23339	54	all	uniqlgene	uniqltag100	100
all.uniqlgene.uniqltag200	40	75	8418	15274	8119	all	uniqlgene	uniqltag200	200
all.uniqlgene.uniqltag200	45	75	7388	15549	7844	all	uniqlgene	uniqltag200	200
all.uniqlgene.uniqltag200	50	75	5312	16473	6920	all	uniqlgene	uniqltag200	200
all.uniqlgene.uniqltag200	55	75	5516	18000	5393	all	uniqlgene	uniqltag200	200
all.uniqlgene.uniqltag200	60	75	3428	18413	4980	all	uniqlgene	uniqltag200	200
all.uniqlgene.uniqltag200	65	75	1541	20012	3381	all	uniqlgene	uniqltag200	200
all.uniqlgene.uniqltag200	70	75	790	22497	896	all	uniqlgene	uniqltag200	200
all.uniqlgene.uniqltag200	74	75	134	23259	134	all	uniqlgene	uniqltag200	200