

463 **SUPPLEMENTARY MATERIAL**

464

465

| | Prediction | Validated class | Accuracy |
|-----------------------|-------------------|------------------------|-----------------|
| <i>False positive</i> | Not contamination | Contamination | 0 |
| <i>False negative</i> | Contamination | Not contamination | 10.96 |
| <i>True positive</i> | Contamination | Contamination | 15.07 |
| <i>True negative</i> | Not contamination | Not contamination | 73.97 |

466 **Supplemental Table 1.** Accuracy of model-predicted classes versus real class for 74 loci withheld
 467 from all training further separated by true/false positives/negatives.

468

469

470

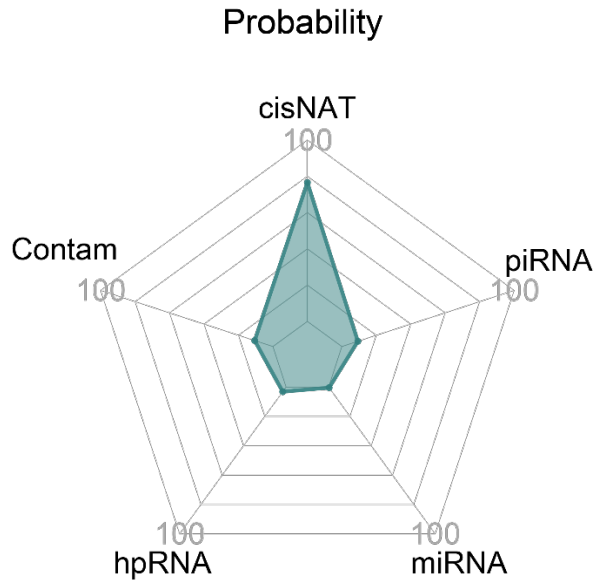
| sRNA class | Accuracy |
|----------------------|-----------------|
| <i>cisNAT</i> | (12/14) 85.71% |
| <i>hpRNA</i> | (9/13) 69.23%* |
| <i>miRNA</i> | (9/11) 81.81% |
| <i>piRNA</i> | (24/24) 100% |
| <i>contamination</i> | (11/11) 100% |

471 **Supplemental Table 2. Model accuracy by sRNA class.** Type I and II error rates of machine
 472 learning predictions by small RNA type for the 74 loci withheld from training.

473

474

475

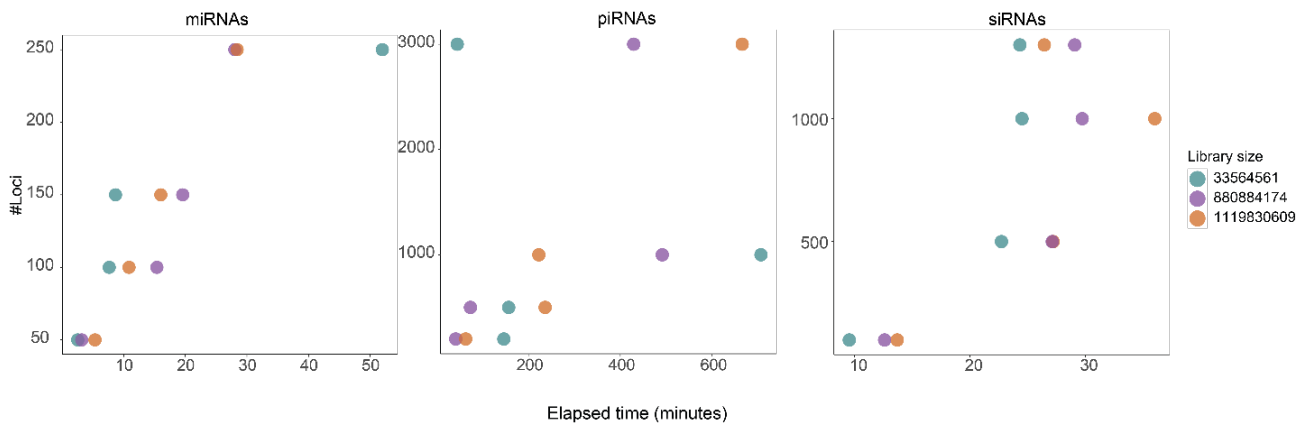


476

477

478 **Supplemental Figure 1. Probability (radar) plot of MiSiPi.RNA machine learning**
 479 **predictions.** Radar plot made using the fmsb package (CRAN).

480



481

482 **Supplemental Figure 2. Program run times (elapsed time) by number of loci and**
 483 **sequencing depth.** Each module of the MiSiPi.RNA package was run on publicly available
 484 datasets of various library sizes. Elapsed time was recorded for each MiSiPi module with basic
 485 arguments in minutes using the system.time() function in R. SRA accessions used were:
 486 SRR22812151 (library size 33564561), SRR23418970 (library size 880884174), and
 487 SRR23418970 concatenated with SRR23418971 before mapping (library size 1119830609).
 488 miRNA and cisNAT siRNA annotations from *Drosophila melanogaster* GCF_000001215.4 were
 489 lifted to the *Drosophila simulans* PacBio genome (Chakraborty et al. 2021) using Liftoff 1.6.2.

490

491 **Supplemental Methods**

492 Accession number of libraries used in the study

| melanogaster | simulans | spiderMite | planarian | zebrafish | mouse |
|--------------------|-------------|-------------|-------------|------------|-------------|
| SRR11680898 | SRR22812151 | SRR10513014 | SRR12426750 | SRR1554490 | SRR13848060 |
| SRR11680899 | SRR22812152 | SRR10513013 | SRR12426752 | SRR1554491 | SRR13848061 |
| SRR11680900 | SRR22812153 | SRR10513015 | SRR12426754 | SRR1554493 | SRR13848062 |
| SRR12313524 | SRR22812154 | SRR10513017 | SRR12426756 | SRR1554494 | SRR24821465 |
| SRR12313525 | SRR22812155 | SRR10513018 | SRR12426758 | | SRR24821466 |
| SRR12313526 | SRR23418970 | SRR10513019 | SRR12426760 | | SRR24821469 |
| | SRR23418971 | | SRR12426751 | | SRR24821470 |
| | | | SRR12426753 | | |
| | | | SRR12426755 | | |
| | | | SRR12426757 | | |
| | | | SRR12426759 | | |
| | | | SRR12426761 | | |
| | | | | | |

493

494 Method for identifying small RNA expressing regions of interest:

495 miRNAs & siRNAs: get reads of only miRNA and siRNA length

```
496 awk 'BEGIN {OFS = "\n"} {header = $0; getline seq; getline qheader ; getline
497 qseq ; if (length(seq) >= 19 && length(seq) <= 23) {print header, seq,
498 qheader, qseq}}' \
499 < trimmed.fq > small.fastq
```

500

501 piRNAs: get reads of piRNA length

502

```
503 awk 'BEGIN {OFS = "\n"} {header = $0; getline seq; getline qheader ; getline
504 qseq ; if(length(seq) >= 23 && length(seq) <= 30) {print header, seq, qheader,
505 qseq}}' \
506 < trimmed.fq > large.fastq
```

507

508 Realign reads to genome

509

510 # all small RNAs

```
511 awk 'BEGIN {OFS = "\n"} {header = $0; getline seq; getline qheader ; getline
512 qseq ; if (length(seq) >= 18 && length(seq) <= 32) {print header, seq,
513 qheader, qseq}}' < trimmed.fq > all.fastq
```

```

514
515 #Align reads to genome
516 bowtie -p 10 -a -m 100 --best --strata --no-unal genome.fna all.fastq -S |
517 samtools view -@ 10 -q 10 -b |samtools sort -@ 10 -m 6G > all.bam
518 samtools index all.bam
519
520 bowtie -p 10 -a -m 100 --best --strata --no-unal genome.fna small.fastq -S
521 | samtools view -@ 10 -q 10 -b |samtools sort -@ 10 -m 6G > small.bam
522 samtools index small.bam
523
524 bowtie -p 20 -a -m 100 --no-unal genome.fna large.fastq -S | samtools view
525 -@ 10 -q 10 -b | samtools sort -@ 10 -m 6G > large.bam
526 samtools index large.bam
527
528 Calculate coverage of all reads over genome
529
530 bedtools genomecov -bg -ibam all.bam | awk '$4 > 100' > HE.tmp.bedgraph
531
532 Get large regions of high expression
533
534 bedtools merge -d 500 -i HE.tmp.bedgraph > HE.tmp.merge.bed
535
536 awk '{n=$2; x=$3; print $1"\t"$2"\t"$3"\t"x-n}' < HE.tmp.merge.bed |
537 awk '$4 > 40' > HE.all.bed #all highly expressed
538
539
540
541 Get potential regions of high miRNA/siRNA RNA expression
542
543 bedtools multicov -bams small.bam all.bam -bed HE.all.bed |
544 awk '{n=$5; x=$6; print $1"\t"$2"\t"$3"\t"n/x}' |
545 awk '$4 > 0.5' > HE.small.bed # only high expressed 21-23 nt
546
547 Get potential regions of high piRNA expression
548
549 bedtools multicov -bams large.bam all.bam -bed HE.all.bed |
550 awk '{n=$5; x=$6; print $1"\t"$2"\t"$3"\t"n/x}' | awk '$4 > 0.5' >
551 HE.large.bed #only high expressed 23-30 nt
552
553 Running the MiSiPi.RNA package:
554
555 Install MiSiPi.RNA. First install devtools and BiocManager
556
557 install.packages("devtools")
558
559 if (!require("BiocManager", quietly = TRUE))
560     install.packages("BiocManager")
561
562 devtools::install_github("stupornova33/MiSiPi.RNA")
563

```

```

564 Running the MiSiPi.RNA on select loci (ROI) and alignments (bam file
565
566     library(MiSiPi.RNA)
567
568     vars <- set_vars(
569         roi = "path/to/bed", # A bed file listing your regions of interest
570         bam_file = "path/to/bam", # A BAM file of aligned reads. Index file must
571         also be present
572         genome = "path/to/genome", # A genome fasta file. Chromosome names must
573         match the BED and BAM files
574         plot_output = TRUE, # (TRUE or FALSE) If TRUE, MiSiPi.RNA will output
575         plots as pdfs
576         path_to_RNAfold = "path/to/ViennaRNA/RNAfold", # Full path to RNAfold
577         executable
578         path_to_RNAplot = "path/to/ViennaRNA/RNAplot", # Full path to RNAplot
579         executable
580         pi_pal = "BlYel", # Palette option for the generated piRNA heatmap (see
581         below)
582         si_pal = "RdYlBl", # Palette option for the generated siRNA heatmap (see
583         below)
584         annotate_region = TRUE, # (TRUE or FALSE) Plots annotated gene features
585         below the hairpin arc plot which is useful for characterizing cisNAT loci
586         weight_reads = "none", # Determines if read counts will be weighted.
587         ("none", "locus_norm", "weight_by_prop", or an integer). This functionality
588         is in beta development.
589         gtf_file = "path/to/gtf", # Full path to a 9 column GTF file. Required
590         only if annotate_region is TRUE
591         write_fastas = FALSE, # (TRUE or FALSE) If TRUE, MiSiPi.RNA will
592         write read pairs from functions to file. Default is FALSE
593         out_type = "pdf" # ("pdf" or "png"),
594         use_bed_names = F, # OPTIONAL (default is FALSE). Specifies whether
595         column 4 of the ROI file contains the locus name to use as file names
596         density_timeout = 3600 # OPTIONAL (default is 3600s), for functions that
597         can run for a long time, move on to next locus after this amount of time
598
599     )
600
601     # Palettes:
602     # Palette options are:
603
604     # "RdYlBl", "BlYel", "yelOrRed", "MagYel", "Greens"
605
606
607     # Main Usage:
608     # choose one method or default is "all"
609     misipi_rna(vars, method = c("miRNA", "piRNA", "siRNA"))
610
611 Using ML module to evaluate loci identity
612
613     # After running with method = "all"

```

```

614 # To make the table for input to the machine learning model, run:
615 # The "_ml.txt" will be found in the misipi_output/date_timestamp/ directory.
616 ml_probability(path_to_table = "misipi_output/date_timestamp/", table =
617 "table_ml.txt")
618
619 # After any method, to create a sortable html summary with plots:
620 make_html_summary(path_to_tables = " misipi_output/date_timestamp/", type =
621 c("siRNA", "piRNA" or "miRNA"), ml_plots = FALSE) # if you didn't run the
622 ml_probability model
623 make_html_summary(path_to_tables = " misipi_output/date_timestamp/", type =
624 c("siRNA", "piRNA" or "miRNA"), ml_plots = TRUE) # if you ran the
625 ml_probability model
626
627 # Optional:
628 # For converting the .ps or .eps output files from the miRNA module to .png,
629 install ImageMagick and ghostscript, then run
630
631 ps2png(path_to_magick_exe, file_dir) or
632 eps2png(path_to_magick_exe, file_dir)
633

```

Changing the ML model for custom analysis

```

634
635 library(dplyr)
636
637
638 # All commented out code is optional and used for troubleshooting accuracy
639 or exploring your data.
640 # Read in output from misipi_rna(vars, method = "all") plus a column
641 containing the class and a column containing the total # of reads in the
642 bam file.
643 # This will be used for library size normalization of read counts.
644 # This file is in "run_all/" and will have the suffix "_ml.txt"
645 # and a column containing the library size of the BAM file used to make the
646 table
647 tbl <- read.table("path/to/machine/learning/output/your_data_ml.txt", sep =
648 "\t", header = TRUE)
649
650 # make sure rows are unique
651 tbl <- tbl %>% dplyr::distinct(log_shap_p, auc, strand_bias, perc_GC,
652 ave_size, perc_first_nucT, perc_A10,
653 highest_si_col, si_dicerz, num_si_dicer_reads,
654 hp_perc_paired, hp_phasedz, hp_mfe,
655 hp_dicerz, mi_perc_paired, mirna_dicerz, mirna_mfe,
656 mirna_overlapz, pingpong_col, max_pi_count,
657 max_piz_overlap, pi_phasedz, pi_phased2z,
658 unique_read_bias, type, species, libsize, X,
659 .keep_all = TRUE)
660
661 # Deal with missing max_piz_overlaps
662 idx <- which(is.na(tbl$max_piz_overlap))
663 tbl$max_piz_overlap[idx] <- -33

```

```

664
665   idx <- which(is.na(tbl$hp_mfe))
666   idx2 <- which(tbl$hp_mfe == 0)
667
668   #set missing hp MFE's to unlikely number
669   tbl$hp_mfe[idx] <- 10000
670   tbl$hp_mfe[idx2] <- 10000
671
672   idx <- which(tbl$mirna_mfe == 0)
673   tbl$mirna_mfe[idx] <- 1000
674
675   for(i in 1:nrow(tbl)){
676     if(tbl$max_pi_count[i] == -33){
677       tbl$max_pi_count[i] <- 0
678     }
679   }
680
681   # Group all types of contamination into one class
682   for(i in 1:nrow(tbl)){
683     if(tbl$type[i] == "snoRNA" || tbl$type[i] == "tRNA" || tbl$type[i] ==
684 "rRNA" || tbl$type[i] == "sno"){
685       tbl$type[i] <- "contamination"
686     }
687   }
688
689   # Normalize counts by libsize
690
691   tbl <- tbl %>% dplyr::mutate(pi_norm = (max_pi_count/libsize/1000000),
692   si_norm = (num_si_dicer_reads/libsize/1000000))
693
694   all_tbl <- tbl %>% dplyr::select(-c(locus_length, species, X, max_pi_count,
695   num_si_dicer_reads, libsize))
696   all_tbl$pi_norm <- as.numeric(all_tbl$pi_norm)
697   all_tbl$si_norm <- as.numeric(all_tbl$si_norm)
698
699   #write.table(all_tbl, file = "new_ML_contam_more_libsize_norm_mod.txt", sep
700 = "\t", quote = FALSE, row.names = FALSE, col.names = TRUE)
701   # install required packages if running for the first time
702   library(dplyr) # for data cleaning
703   library(CatEncoders)
704   library(xgboost)
705   library(e1071)
706
707   # The MFE doesn't need to be negative, so take absolute value
708   all_tbl$mirna_mfe <- abs(all_tbl$mirna_mfe)
709
710   all_tbl$hp_mfe <- abs(all_tbl$hp_mfe)
711
712   # remove the locus name and
713   # Hot encode the type class

```

```

714 tbl <- all_tbl %>% dplyr::select(-c(locus))
715 train_labs <- LabelEncoder.fit(tbl$type)
716 tbl$type <- transform(train_labs, tbl$type)
717 # transform to be zero-based
718 tbl$type <- tbl$type - 1
719
720 # check whether there is a class imbalance
721 table(tbl$type) # Raw class counts
722 prop.table(table(tbl$type)) * 100 #
723
724 # If serious class imbalance exists, consider methods to either re-sample,
725 # add more data, or use class weighting.
726 # tbl$type <- as.factor(tbl$type) # caret requires a factor
727
728 # or can use SMOTE:
729 # Better than simple upsampling – creates synthetic examples of minority
730 classes.
731 # remotes::install_github("cran/DMwR")
732 # library(DMwR)
733 #
734 # tbl$type <- as.factor(tbl$type)
735 # balanced_data <- SMOTE(type ~ ., data = tbl, perc.over = 100, perc.under
736 = 200)
737 # tbl <- balanced_data
738
739 # Create Training and Test data -
740 set.seed(100) # setting seed to reproduce results of random sampling
741 # Segregate 80% of data for training and leave 20% out for cross validation
742 trainingRowIndex <- sample(1:nrow(tbl), 0.8*nrow(tbl))# row indices for
743 training data
744 trainingData <- tbl[trainingRowIndex, ] #>% dplyr::select(-c(locus,
745 species)) # model training data
746 testData <- tbl[-trainingRowIndex, ] #>% dplyr::select(-c(locus,
747 species)) # test data
748
749 train_x <- model.matrix(type ~ ., trainingData)[, -1]
750
751 test_x <- model.matrix(type ~ ., testData)[, -1]
752
753 train_y <- trainingData$type
754
755 test_y <- testData$type
756
757 train_y <- as.integer(as.character(train_y)) # safely convert factor to
758 numeric
759 test_y <- as.integer(as.character(test_y))
760 # convert the train and test data into xgboost matrix type.
761 dtrain = xgb.DMatrix(data=train_x, label=train_y)
762 dval = xgb.DMatrix(data=test_x, label=test_y)
763

```

```

764 # Store X and Y for later use
765 #y = train_y
766 #x = train_x
767
768 options(scipen = 999)
769
770 #featurePlot(x = x,
771 #           y = as.factor(y),
772 #           plot = "density",
773 #           strip = strip.custom(par.strip.text = list(cex = .7)),
774 #           scales = list(x = list(relation="free"),
775 #                          y = list(relation="free")))
776
777 ## read in the cross-validation table and transform the values same as the
778 training data
779 # I made my own cross-validation with completely unseen data to get a
780 better idea of real life accuracy,
781 # as the held-back training data seemed to have inflated accuracies.
782 # This table is formatted the same as the ml table.
783
784 Dat <-
785 read.table("C:/Users/tmjar/Desktop/MiSiPi_Stuff/machine_learning_species/ne
786 w_all_species/table_versions/all_species_w_whitefly_CV.txt", header = TRUE)
787
788 # Group all types of contamination as one class
789 dat <- dat %>% dplyr::select(-c(locus, locus_length, X,species))
790 for(i in 1:nrow(dat)){
791   if(dat$type[i] == "snoRNA" || dat$type[i] == "tRNA" || dat$type[i] ==
792 "rRNA" || dat$type[i] == "sno" || dat$type[i] == "contamination"){
793     dat$type[i] <- "contamination"
794   }
795 }
796
797
798 dat <- dat %>% dplyr::mutate(pi_norm = "", si_norm = "")
799
800 #normalize counts by libsize
801
802 for(i in 1:nrow(dat)){
803   #idx <- which(df$dataset == dat$X[i])
804   dat$pi_norm[i] <- dat$max_pi_count[i]/dat$libsize[i]/1000000
805   dat$si_norm[i] <- dat$num_si_dicer_reads[i]/dat$libsize[i]/1000000
806
807 }
808 dat$pi_norm <- as.numeric(dat$pi_norm)
809 dat$si_norm <- as.numeric(dat$si_norm)
810 dat <- dat %>% dplyr::select(-c(libsize, max_pi_count, num_si_dicer_reads))
811
812 test_labs <- LabelEncoder.fit(dat$type)
813 dat$type <- transform(test_labs, dat$type)

```

```

814 dat$type <- dat$type - 1
815
816 idx <- which(is.na(dat$max_piz_overlap))
817 dat$max_piz_overlap[idx] <- -33
818
819 # Set missing MFE to unlikely number
820
821 idx1 <- which(is.na(dat$hp_mfe))
822 idx2 <- which(dat$hp_mfe == 0)
823 dat$hp_mfe[idx1] <- 10000
824 dat$hp_mfe[idx2] <- 10000
825
826 idx1 <- which(is.na(dat$mirna_mfe))
827 idx2 <- which(dat$mirna_mfe == 0)
828
829 dat$mirna_mfe[idx1] <- 1000
830 dat$mirna_mfe[idx2] <- 1000
831
832 dat$hp_mfe <- abs(dat$hp_mfe)
833
834 new_data_y <- dat$type
835 #new_data_y <- new_data_y - 1
836 new_data_x <- dat %>% dplyr::select(-type)
837
838 ##### Do the training
839
840 rounds <- c(700,800,900,1000)
841 etas1 <- c(seq(0.01, 0.3, 0.05))
842 #etas <- c(etas1, etas2)
843 #etas <- c(seq(0.01, 0.4, 0.1))
844 subs <- c(0.7,0.75,0.8, 0.85,0.9, 0.95, 0.1)
845 depths <- c(5,6,7,8,9,10,11)
846 lambdas <- c(0.01,0.02,0.03,0.04,0.05,0.06,0.07,0.08,0.09,0.1,1,0)
847 alphas <- c(0,0.1,0.5,1.0)
848 gammas <- c(0.05,0.06,0.07,0.08,0.09,0.1,0.3,0.5,0.7,0.9,1.0)
849
850 len <- length(gammas)
851 n_iter <- len
852
853 pb <- txtProgressBar(min = 0,
854                     max = n_iter,
855                     style = 3,
856                     width = 50,
857                     char = "=")
858
859 results <- data.frame(par = numeric(length = length(len)), mean_accuracy =
860 numeric(length = length(len)))
861
862 for( i in 1:length(gammas)){
863   set.seed(1234)

```

```

864     #run with softmax for quick accuracy check, pick best model, re-run with
865     softprob and save
866     cur_model <- xgboost(data = dtrain, max.depth = 6, subsample = 0.85, eta
867     = 0.01, nthread = 2, lambda = 0.01, alpha = 0.5,
868     gamma = 0.05, eval_metric = "mlogloss", nrounds =
869     700, objective = "multi:softmax", num_class = 5, verbose = 0,
870     early_stopping_rounds = 10, print_every_n = 10)
871
872     # once final model is selected, run the following to get a model that
873     outputs percent probability instead of just "class"
874     # prob_model <- xgboost(data = dtrain, max.depth = 6, eval_metric =
875     "mlogloss", eta = 0.01, subsample = 0.85, nthread = 2,
876     # alpha = 0.1, lambda = 0.01, gamma = 0.05,
877     # nrounds = 700, objective = "multi:softprob",
878     num_class = 5, verbose = FALSE)
879
880     # then save prob_model to a .rds file
881     real_pred <- predict(cur_model, as.matrix(new_data_x))
882     #df <- as.data.frame(matrix(real_pred, ncol = 2, byrow = TRUE)) #>%
883     dplyr::select(c(V1))
884
885     prob_mat <- matrix(real_pred, ncol = 5, nrow = nrow(new_data_x))
886     cur_err <- mean(real_pred != new_data_y)
887     cur_accuracy <- 1 - cur_err
888     print(cur_accuracy)
889
890     results[i,] <- c(gammas[i], cur_accuracy)
891     setTxtProgressBar(pb, i)
892 }
893
894 # To check whether model is overfit
895 # log <- model$evaluation_log
896 # plot(log$iter, log$train_mlogloss, type = "l", col = "blue")
897 # plot(log$iter, log$eval_mlogloss, type = "l", col = "red")
898 # legend("topright", legend = c("Train", "Validation"), col = c("blue",
899 "red"), lty = 1)
900
901 # save the model
902 xgb.save(prob_model, "path/to/xgb_tuned.rds")
903
904 # load the model
905 all_model <- xgb.load("path/to/xgb_tuned.rds")
906
907 # Calculate Type I and Type II error rates
908 all_pred <- predict(all_model, as.matrix(new_data_x))
909 all_df <- as.data.frame(matrix(all_pred, ncol = 5, byrow = TRUE))
910 #all_pred <- all_pred[,1]
911 colnames(all_df) <- c("prob_cis", "prob_contam", "prob_hp", "prob_mi",
912 "prob_pi")
913 all_df$real_type <- dat$type

```

```
914
915 for(i in 1:nrow(all_df)){
916   real <- all_df$real_type[i]
917   #true and contam
918   #true and not contam
919   #false and contam
920   #false and not contam
921   max <- which(all_df[i,1:5] == max(all_df[i,1:5]))
922   if((max - 1) == real){
923     all_df$cond[i] <- "TRUE"
924   } else {
925     all_df$cond[i] <- "FALSE"
926   }
927
928   if(all_df$real_type[i] == 1){
929     all_df$contam[i] <- "contam"
930   } else {
931     all_df$contam[i] <- "not_contam"
932   }
933 }
934
```