# Supplementary Figures
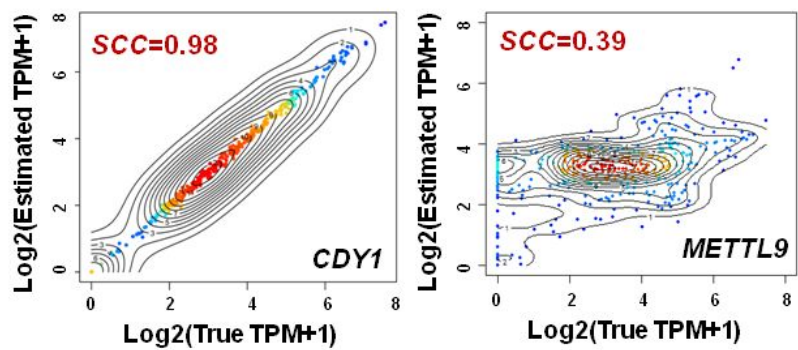
**Systematic assessment of long-read RNA-seq methods for transcript identification and quantification**
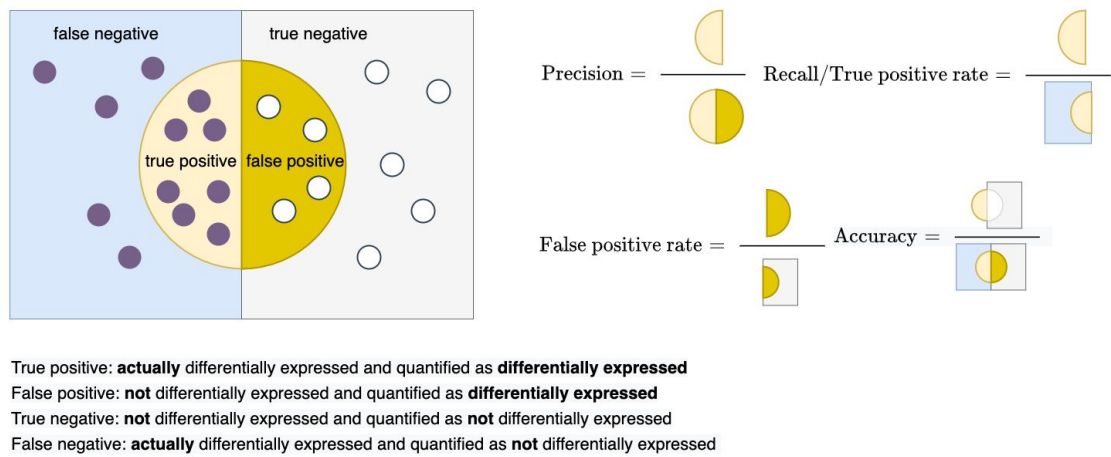
Pardo-Palacios et al.
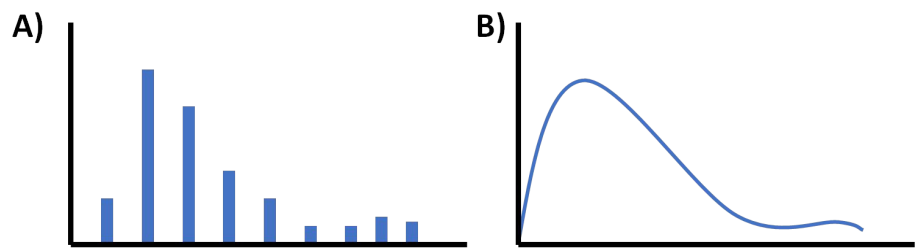
## Spearman Correlation Coefficient



**Supplementary Fig. 1: Spearman Correlation Coefficient (*SCC*) between the estimation and gold standard.** The simulation study based on *SCC* reveals gene *CDY1* can be accurately quantified but not gene *METTL9*.

**Fold change based evaluation**



True positive: **actually** differentially expressed and quantified as **differentially expressed**
False positive: **not** differentially expressed and quantified as **differentially expressed**
True negative: **not** differentially expressed and quantified as **not** differentially expressed
False negative: **actually** differentially expressed and quantified as **not** differentially expressed
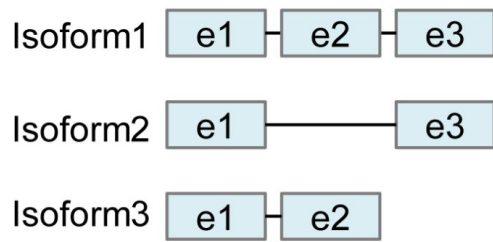
**Supplementary Fig. 2: Log-fold-change-based evaluation metrics.** This figure illustrates how ROC statistics such as precision, recall and accuracy are calculated. It measures the performance on detecting real biological changes.

## Resolution Entropy



**Supplementary Fig 3. Resolution Entropy.** (A) The software output only a few certain discrete values has lower resolution entropy as it cannot capture the continuous and subtle difference of gene expressions. (B) The software with continuous output values has higher resolution entropy.

## Gene Structure

Isoform1 | e1 | e2 | e3 |
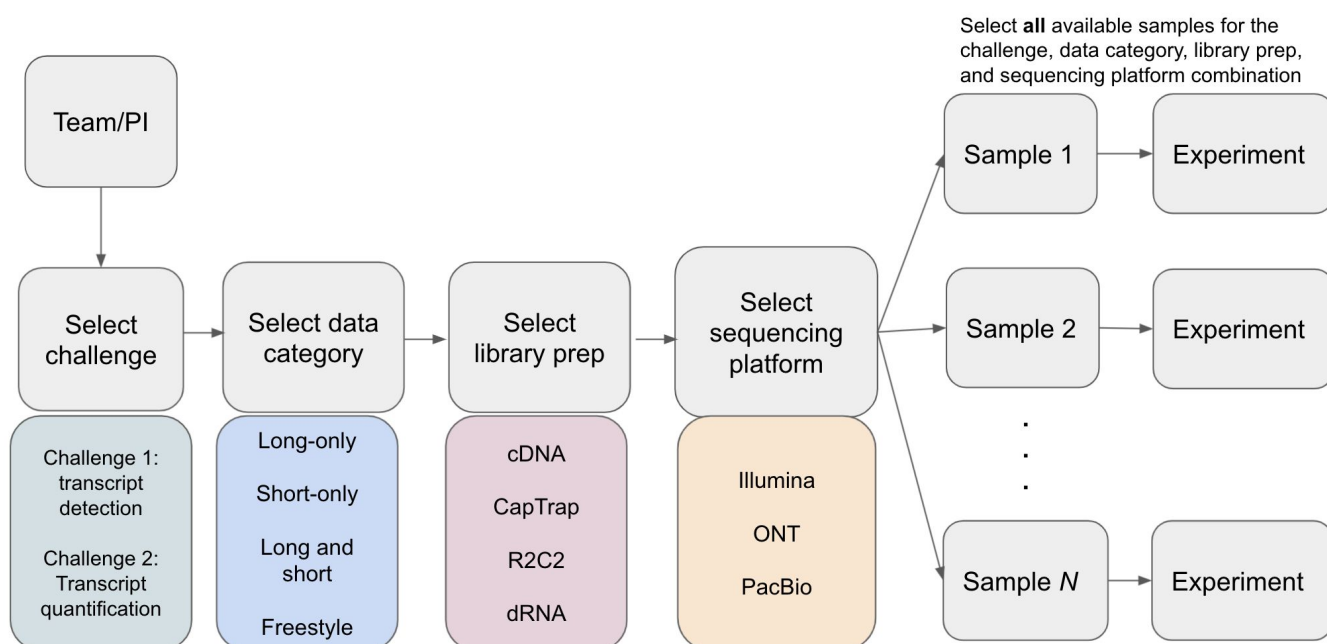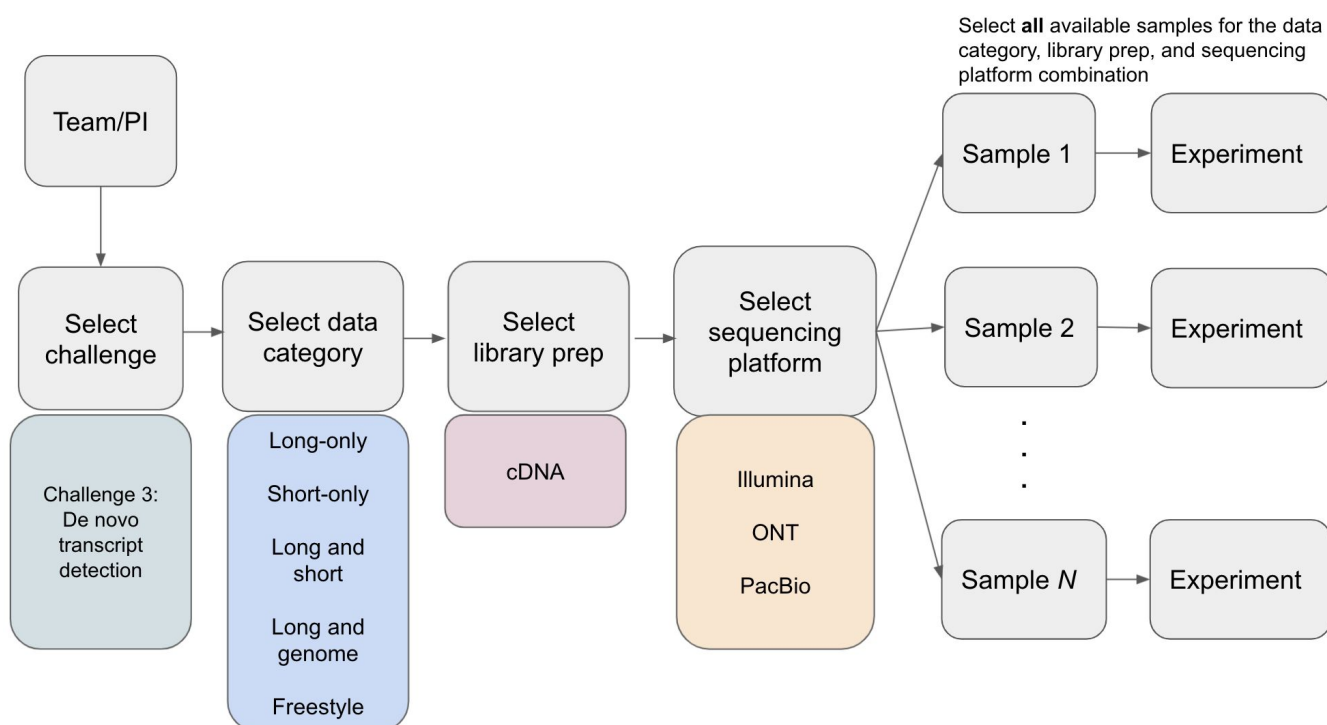
Isoform2 | e1 |————| e3 |

Isoform3 | e1 | e2 |

## Isoform-exon binary matrix A

$$A = \begin{array}{c c} & \begin{array}{ccc} e1 & e2 & e3 \end{array} \\ \begin{array}{c} \text{Isoform1} \\ \text{Isoform2} \\ \text{Isoform3} \end{array} & \begin{bmatrix} 1 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix} \end{array}$$
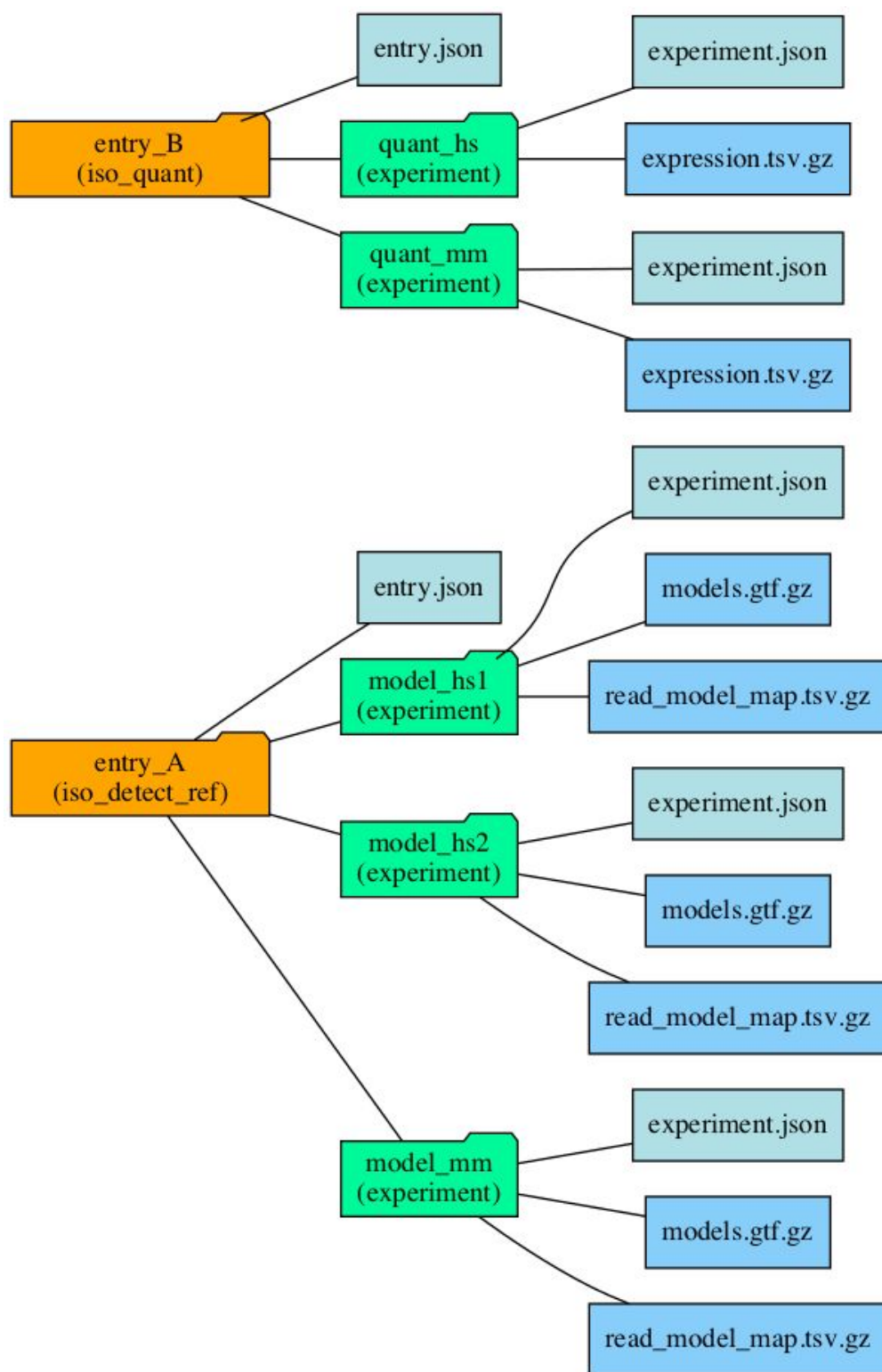
## K-value

$$K(A) = \frac{\sigma_{max} A}{\sigma_{min} A}$$

**Supplementary Fig. 4: Description of K-value.** A measure of the complexity of exon-isoform structures for each gene.
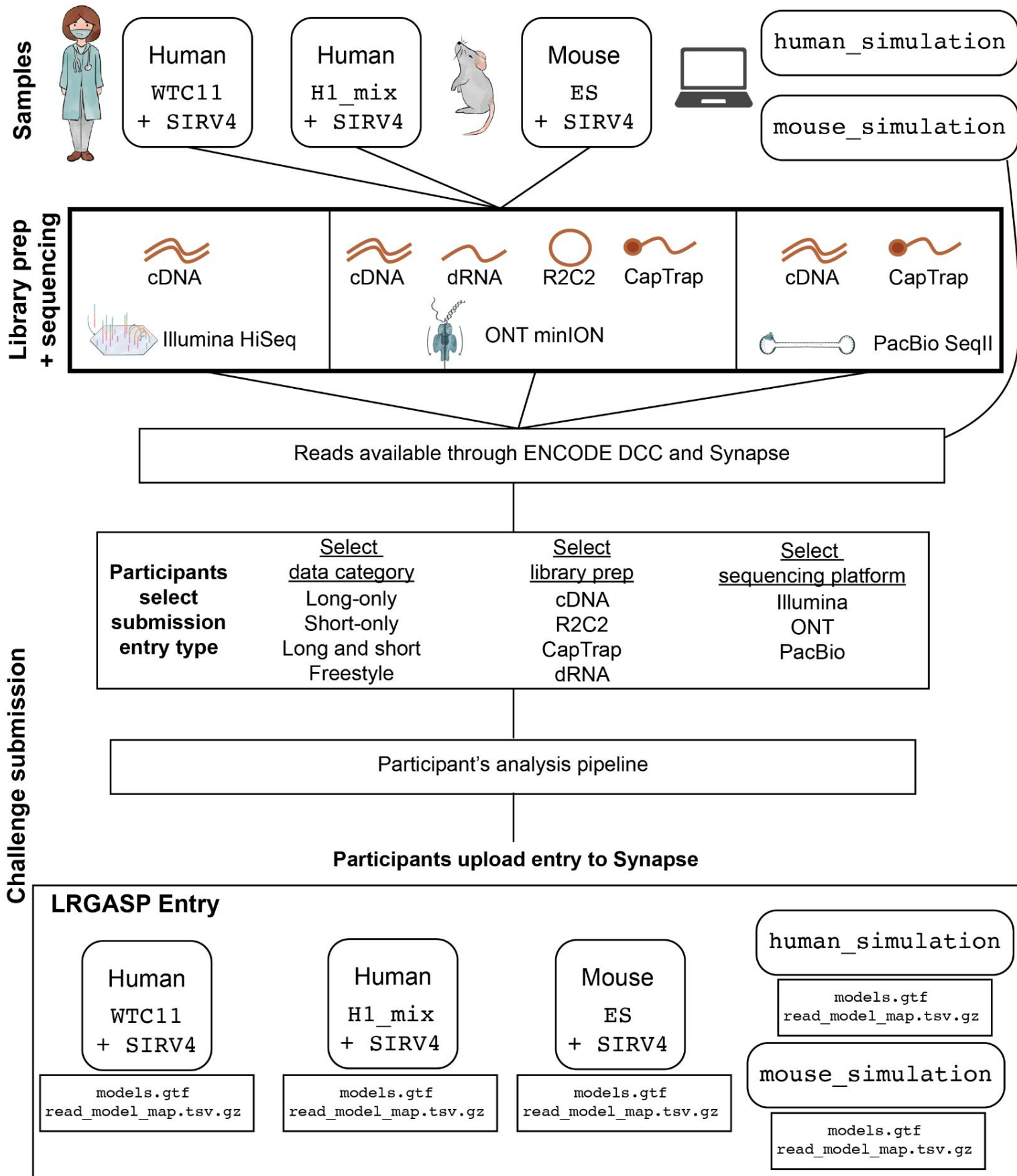
**a**



**b**



**Supplementary Fig. 5: Challenge submission. a,** Overview of submissions to Challenges 1 and 2. Each entry will be derived from a specific data category, library prep, and sequencing platform combination. All available samples for the selected combination must be included in an entry **b**, Overview of submissions for Challenge 3
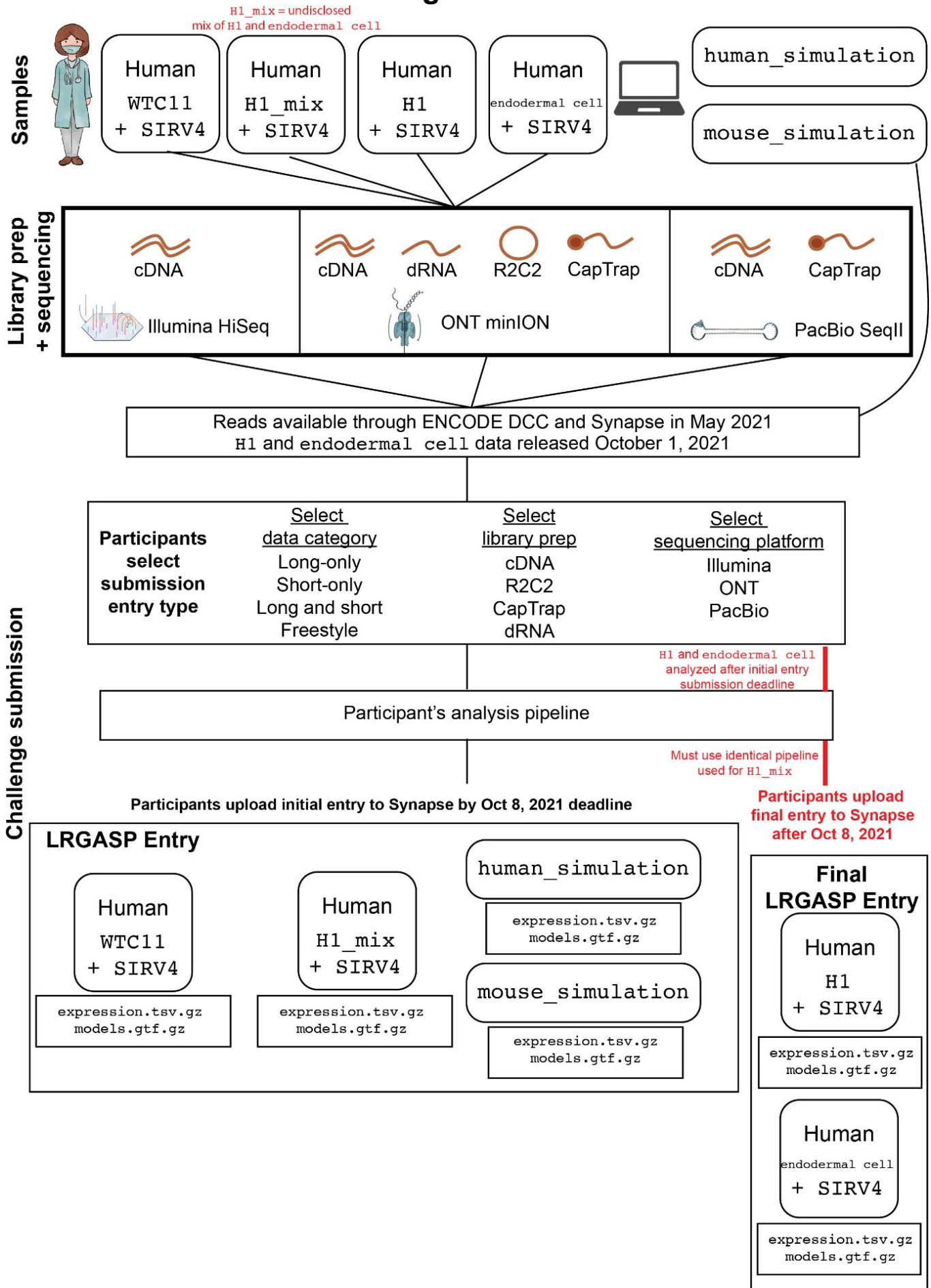
**Supplementary Fig. 6. Schematic of directory structure and files that would be included in each entry**

# Challenge 1 Overview



**Supplementary Fig. 7: Flow diagram of Challenge 1: Transcript isoform detection with a high-quality genome.** Samples, library prep methods, and sequencing platforms used in the challenge are indicated at the top. Participants select which data category, library prep, and sequencing platform to analyze, run their pipelines to generate transcript predictions, and submit an entry which includes predictions for all samples. The entries include a .gtf file of the transcript models and a .tsv file that assigns reads that supported each transcript model.

# Challenge 2 Overview



**Supplementary Fig. 8. Flow diagram of Challenge 2: Transcript isoform quantification.** Samples, library prep methods, and sequencing platforms used in the challenge are indicated at the top. Participants select which data category, library prep, and sequencing platform to analyze, run their pipelines to generate transcript predictions, and submit an entry which includes predictions for all samples. The entries include a .gtf file of the transcript models that are quantified and a .tsv file of the expression quantification. The H1 and endodermal cell samples were released after the initial submission deadline and participants were required to submit the quantification after the deadline.
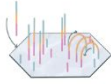
# Challenge 3 Overview



**Supplementary Fig. 9. Flow diagram of Challenge 3.** Samples, library prep methods, and sequencing platforms used in the challenge are indicated at the top. Partici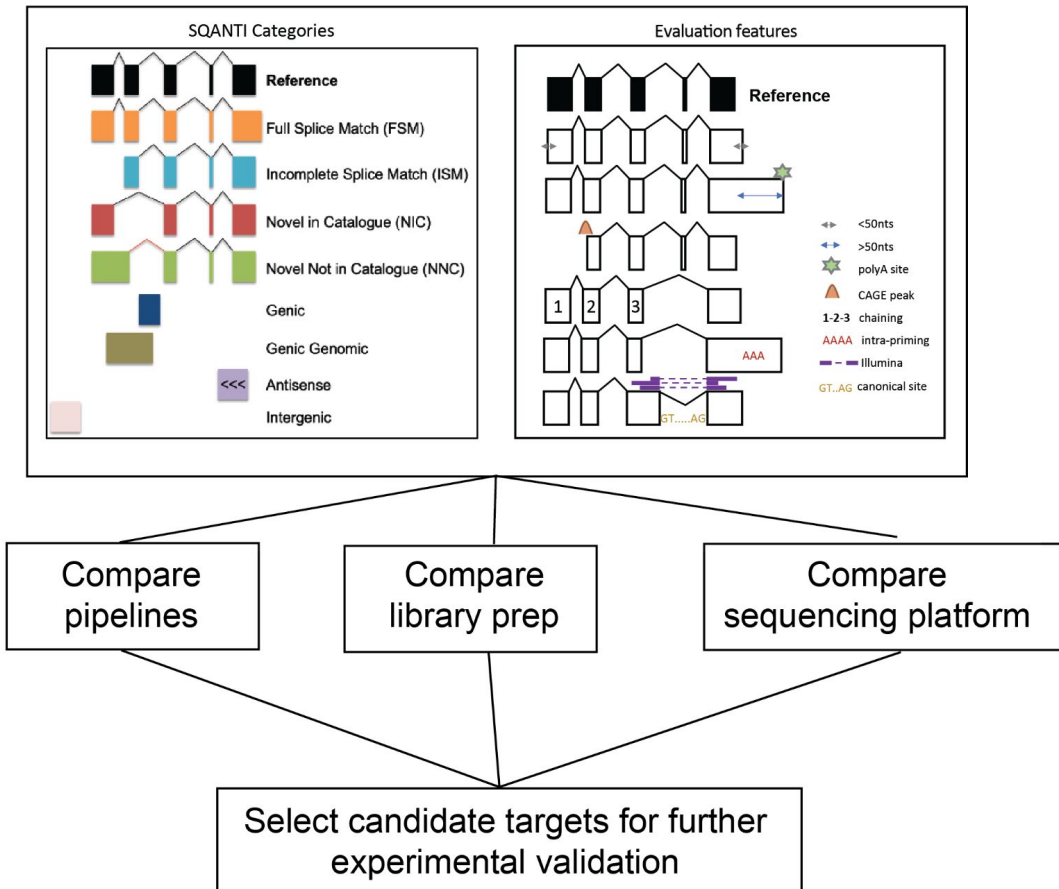pants select which data category and sequencing platform to analyze, run their pipelines to generate transcript predictions, and submit an entry which includes predictions for all samples. The entries include a .fasta file of the transcript models and a .tsv file that assigns reads that supported each transcript model.
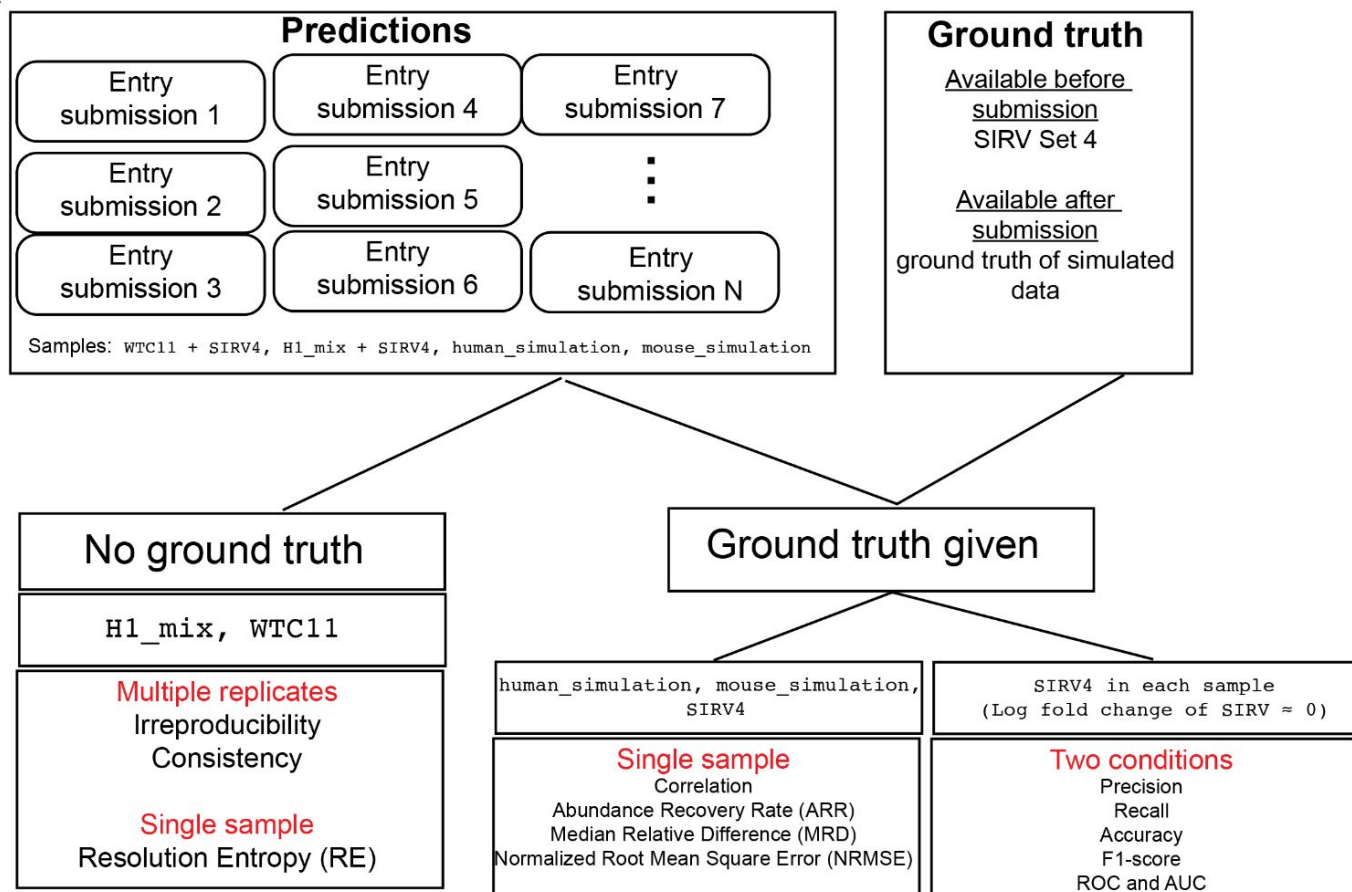
# Challenge 1 Evaluation

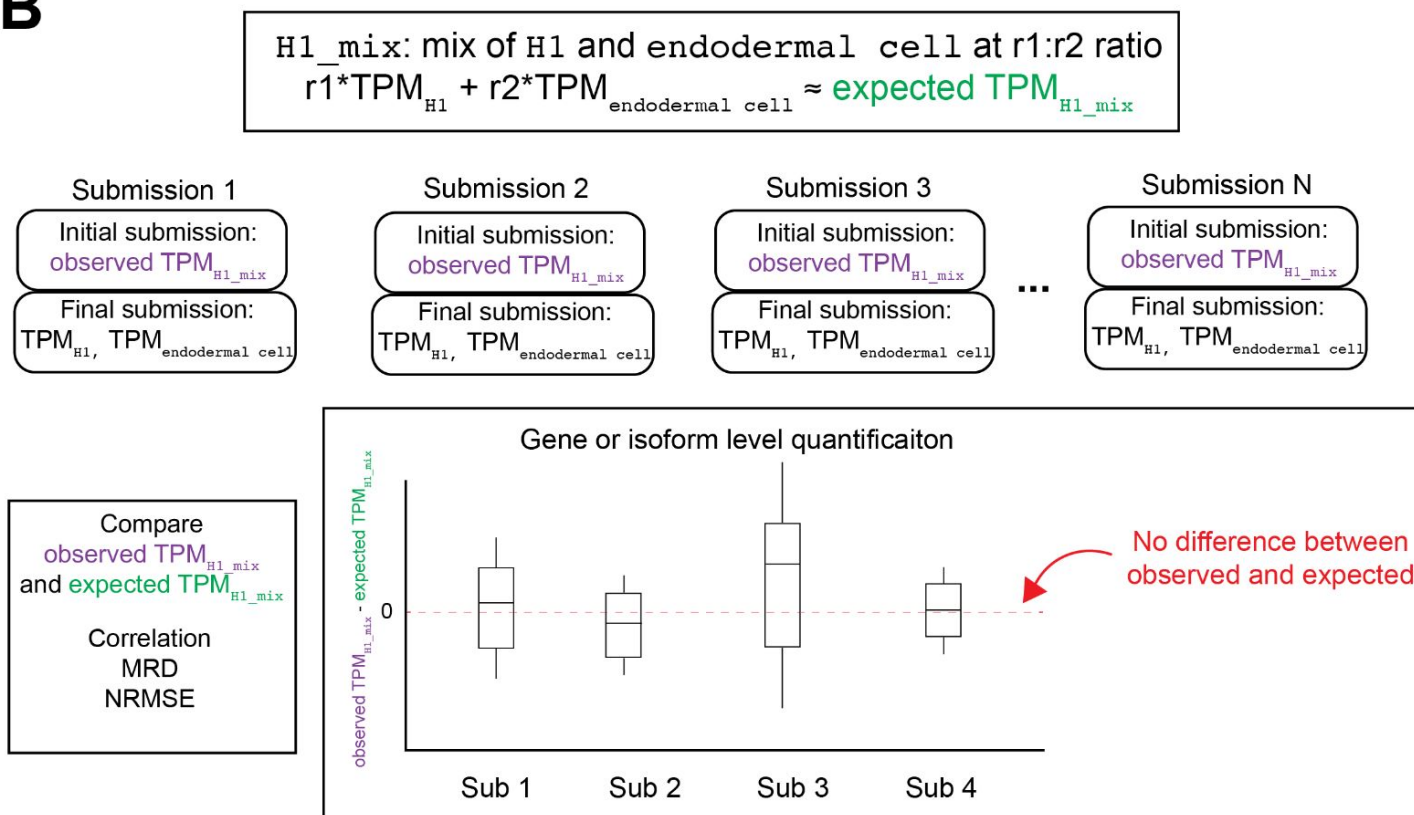

**Supplementary Fig. 10. Flow diagram of the evaluation for Challenge 1.** Benchmarks and additional orthogonal data that will be used for the evaluation are indicated. For example, CAGE and QuantSeq data from WTC11 cells were generated and made available only after participant submissions; therefore, they represent "hidden" data. These will be used to define 5' transcript starts and 3' ends.
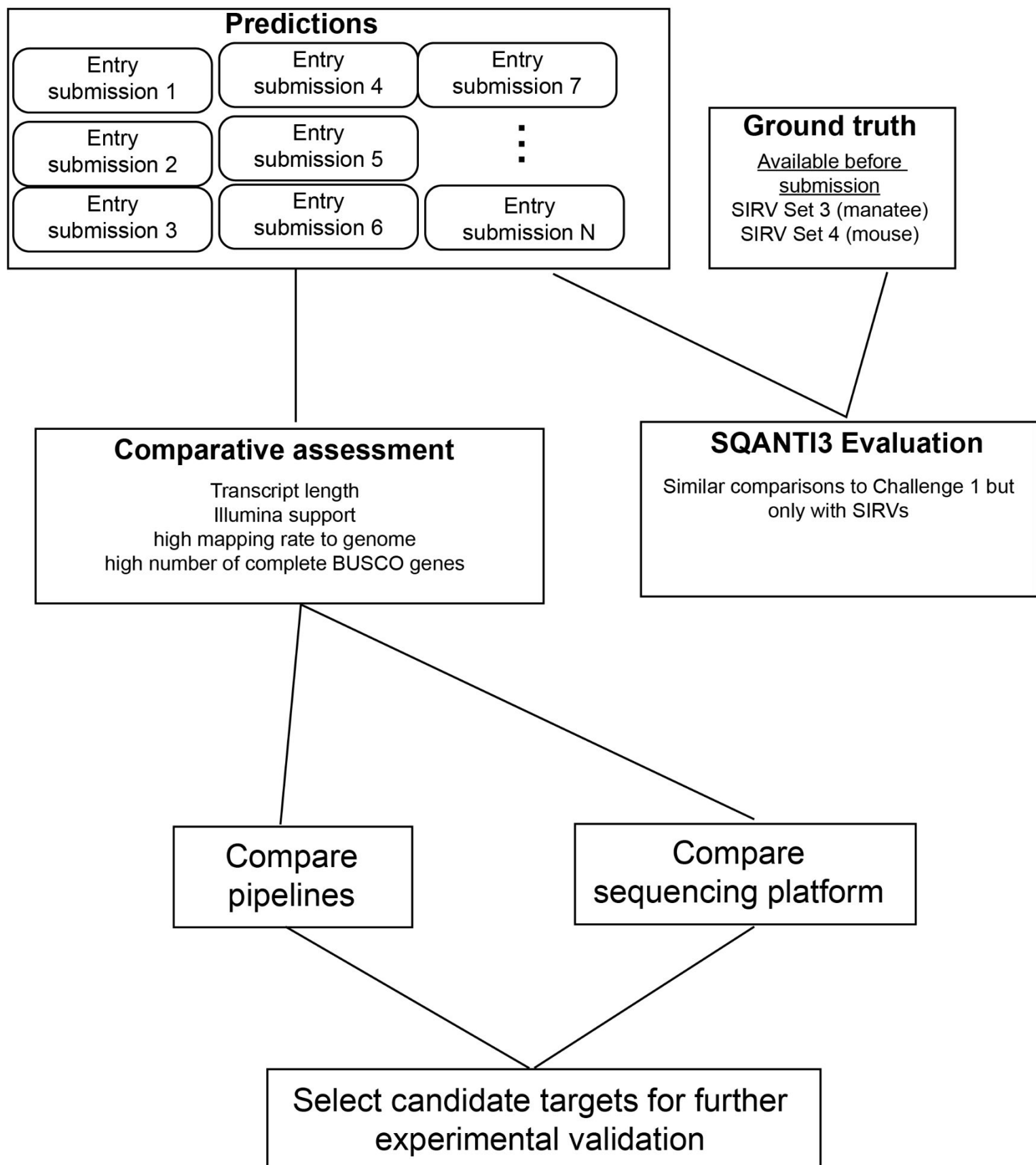
**Supplementary Fig. 11. Flow diagram of the evaluation for Challenge 2. (A)** Evaluation of Challenge 2 can be separated into metrics when a ground truth is known or a ground truth is unknown. (**B**) Example analyses to evaluate transcript expression using the cell mixing experiment. A sample, H1_mix, was initially provided for quantification which was a mix of H1 cells and endodermal cells at an undisclosed ratio. After the initial submission, the individual H1 and endodermal cell samples were released and participants submitted quantifications for each.

# Challenge 3 Evaluation

**Predictions**

| | | |
|---|---|---|
| Entry submission 1 | Entry submission 4 | Entry submission 7 |
| Entry submission 2 | Entry submission 5 | ⋮ |
| Entry submission 3 | Entry submission 6 | Entry submission N |

**Ground truth**

Available before submission
SIRV Set 3 (manatee)
SIRV Set 4 (mouse)

**Comparative assessment**

Transcript length
Illumina support
high mapping rate to genome
high number of complete BUSCO genes

**SQANTI3 Evaluation**

Similar comparisons to Challenge 1 but only with SIRVs

Compare pipelines

Compare sequencing platform

Select candidate targets for further experimental validation

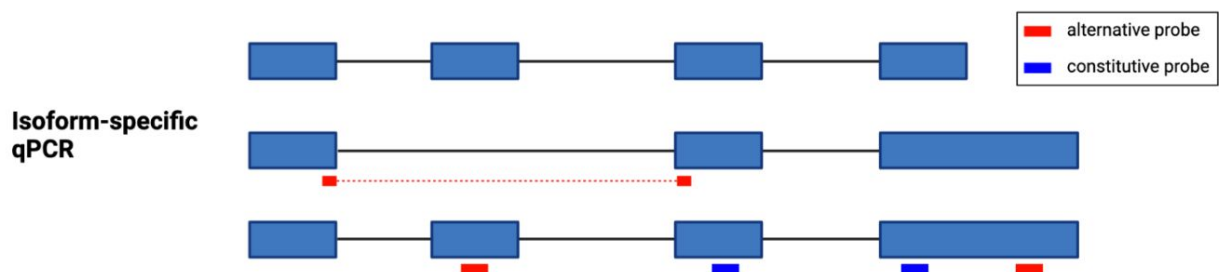**Supplementary Fig. 12. Flow diagram of the evaluation for Challenge 3.** Only SIRVs are available for ground truth information. The evaluation will be based on a comparative assessment of the predictions followed by targeting specific candidates for further validation.
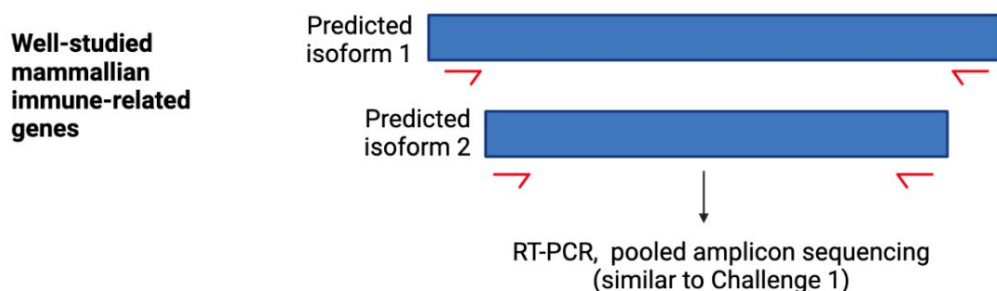
**Supplementary Fig. 13: Experimental validation approaches for the LRGASP challenges. (A)** Multiple categories of types of transcript will be selected for validation (shown in green boxes). These loci will be viewed in the UCSC Genome Browser along with additional datasets to aid in the manual design of primers. Amplicons will be analyzed by fragment size and also pooled to perform long-read sequencing with PacBio and ONT **(B)** A select number of genes will be selected for transcript isoform-specific qPCR. A combination of probes detecting constitutive and alternative regions will be used. (C) RT-PCR validation will be performed similar to Challenge 1, except transcript will be selected from well-studied mammalian immune-related genes.