

Introduction

Long-read RNA sequencing enables direct observation of full-length RNA molecules, revealing transcript diversity not captured by short-read sequence data. We present isocall, an efficient multi-sample isoform detection tool, together with workflows to enable population-scale full-length RNA analysis.

Kinnex full-length RNA

By concatenating transcripts into larger fragment libraries, Kinnex increases throughput for full-length RNA transcript sequencing, enabling large-scale studies at high-resolution. New tools and workflows improve and simplify data analysis.

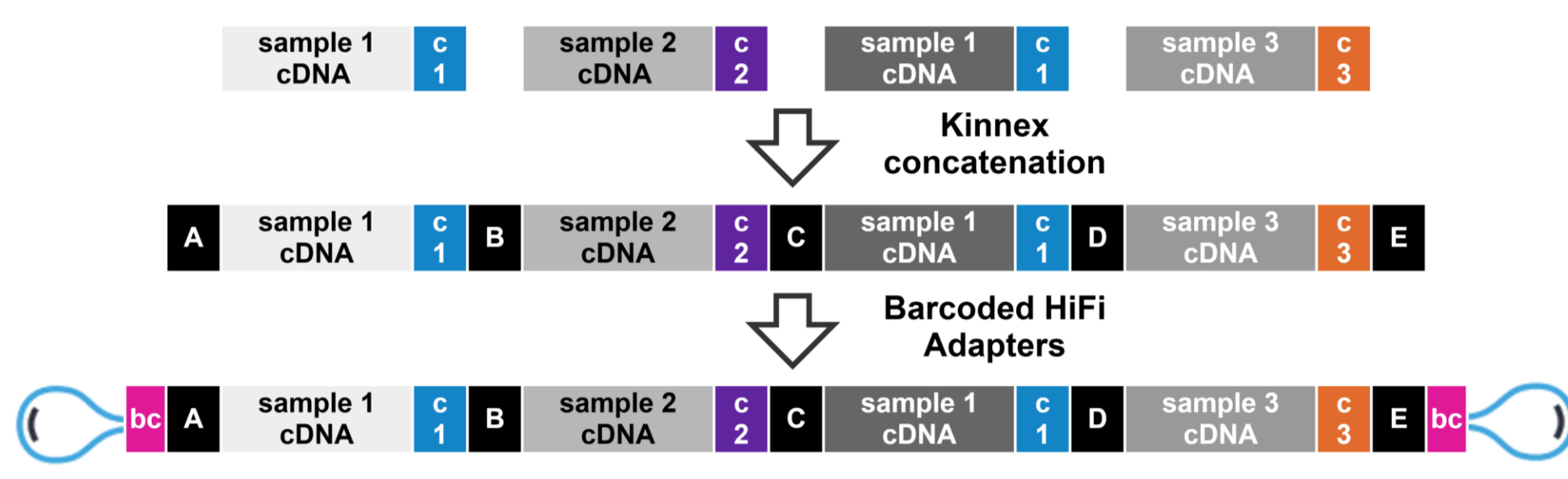


Figure 1. Overview of Kinnex Full-Length RNA. After cDNA generation, Kinnex adapters are attached to the transcripts via PCR. These are then concatenated into a long molecule during array formation. The resulting library can then be cleaned up and sequenced. The library is barcoded at two levels: the original barcode incorporated on the transcript and the HiFi adapter barcode incorporated after concatenation.

Kinnex RNA pre-processing workflow generates transcripts bam for biosamples

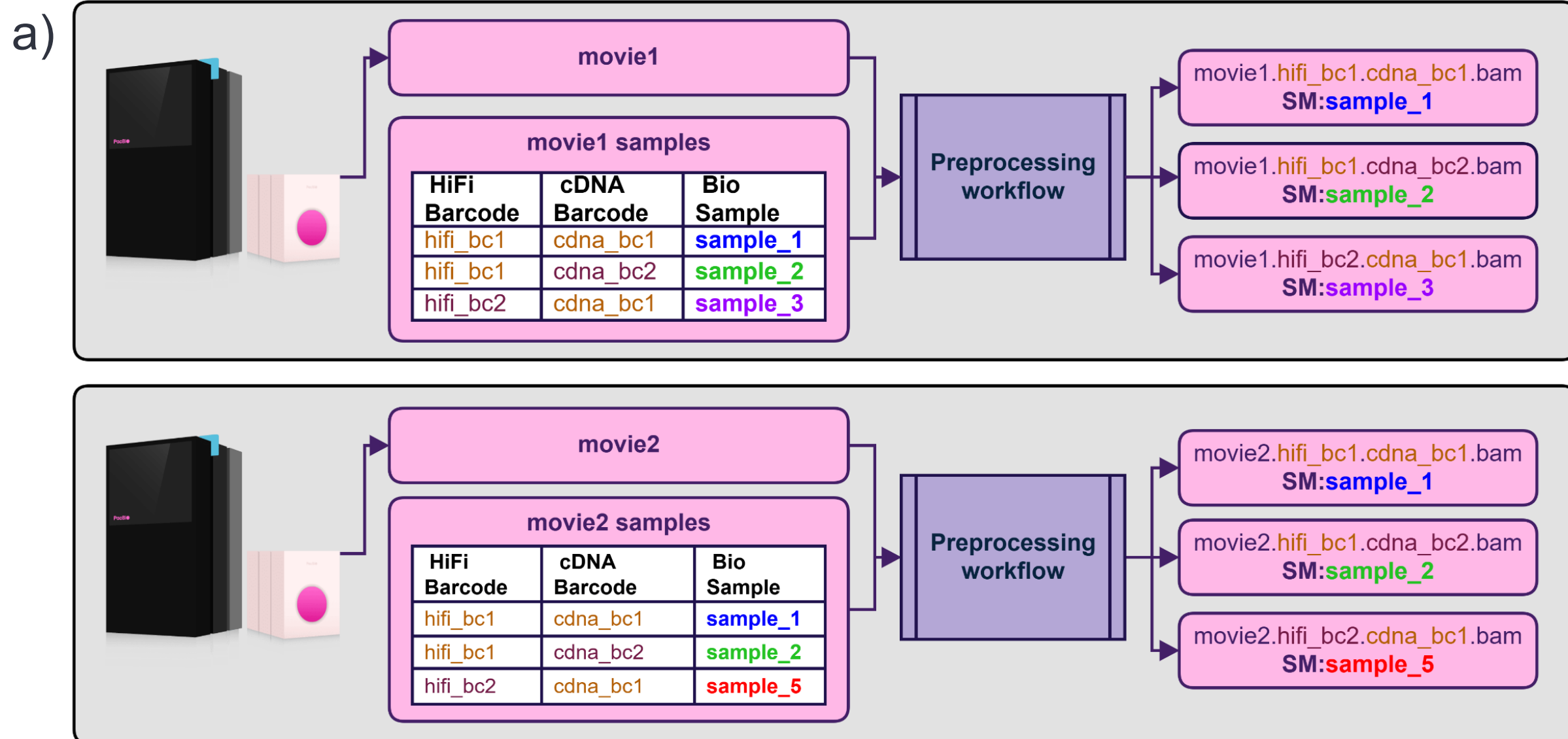


Figure 3: Kinnex RNA preprocessing workflow. The pre-processing workflow allows users to seamlessly obtain transcript BAM files from PacBio long-read RNA sequencing data. a) Users define a two-level biosample CSV specifying both the HiFi and cDNA barcodes for each sample. The workflow generates sample annotated BAM files. b) The preprocessing workflow starts from instrument BAM files. It desegments, demultiplexes, filters concatemers, and trims the poly-A tail to generate per-sample transcript BAM files.

Kinnex RNA secondary analysis workflow generates results for each biosample across acquisitions

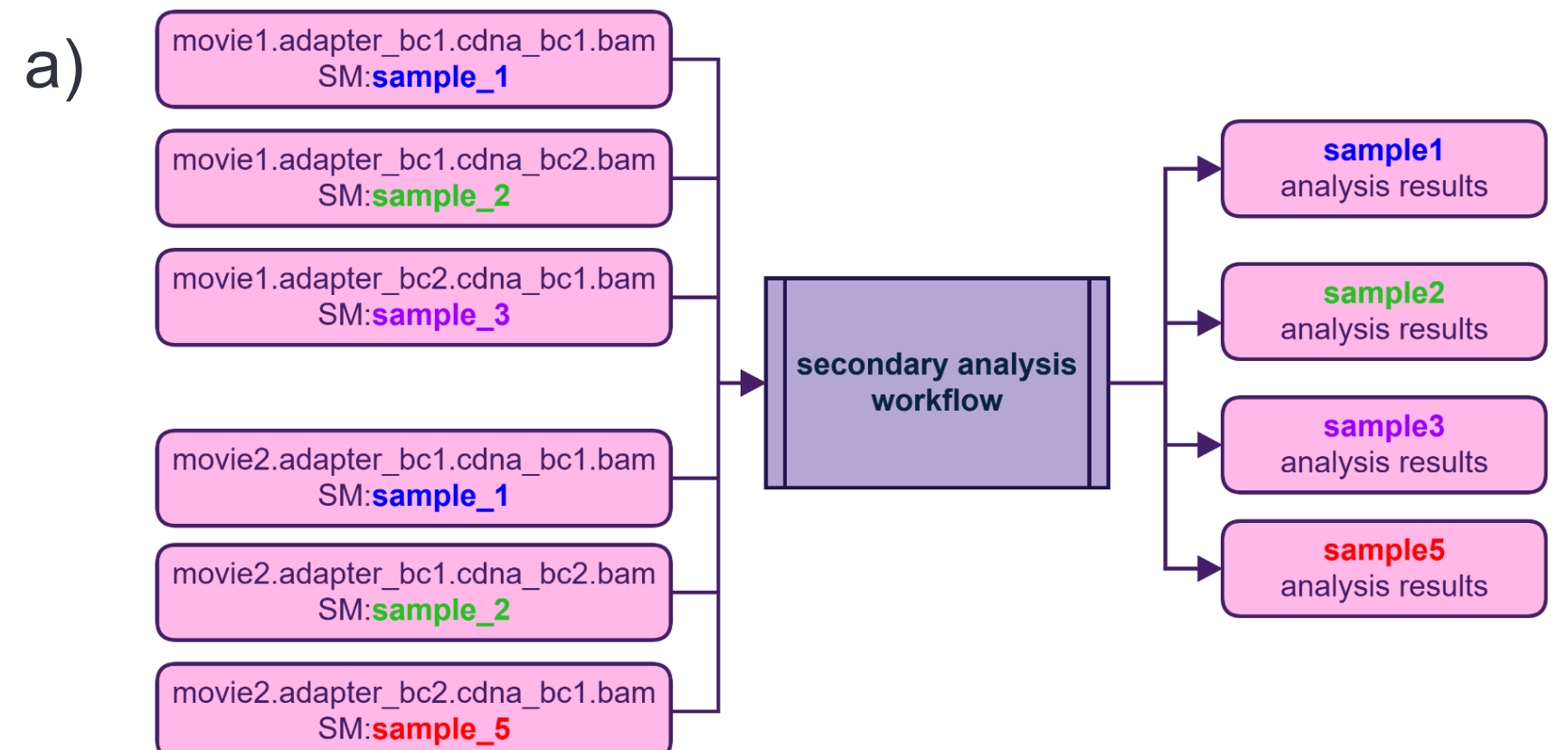


Figure 4. Kinnex RNA Secondary Analysis Workflow. The secondary analysis workflow generates results per biosample regardless of acquisition. a) The secondary analysis workflow takes biosample-annotated transcripts BAMs from any number of acquisitions to generate biosample-level analysis results. In the provided example, sample_1 and sample_2 are sequenced across two acquisitions while sample_3 and sample_5 are only on a single acquisition. b) The secondary analysis workflow first merges the BAM files based on their biosample annotation. Transcript reads are mapped to the reference genome. Joint isoform calling and isoform classification are then performed at the sample level.

Data

We performed isoform joint calling, classification and filtering across two datasets:

- A subset of the 1000 genomes project to assess run time and memory use on large-scale data.
- RNA expression of differentiating stem cell (WTC11) derived primordial endothelial cells from Wissel et. al¹ to look at isocall performance on SIRV data and investigate novel isoforms across time and replicates.

Isocall run time and memory use

Isocall completed isoform joint calling on 250 mapped transcriptomes (~10M S-reads per sample) in ~3 hours using 64 threads and 4 GBs of RAM. Run time and memory use scale linearly.

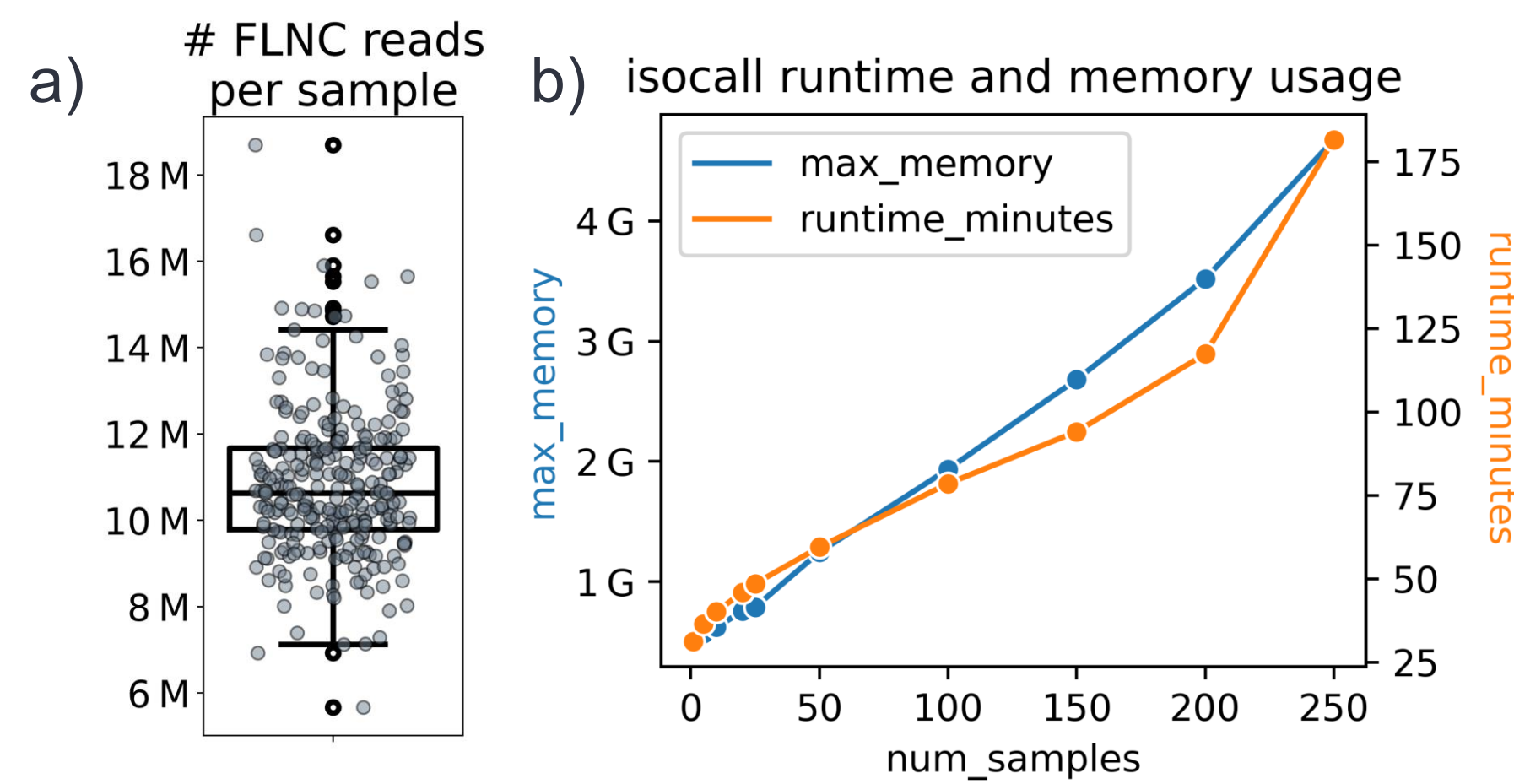


Figure 2. Isocall run time and memory usage on large-scale data. a) samples had ~10M transcript reads each. b) run time and memory on 64 threads, up to 250 samples.

Isocall Performance on SIRV Benchmark

We assessed performance of isocall using the SIRV spike-ins in the iPSC data. Users can control isocall's sensitivity by changing a single parameter which determines the minimum fraction of reads that must be explained by called isoforms at each locus. A higher fraction will yield to a higher sensitivity with a tradeoff on specificity.

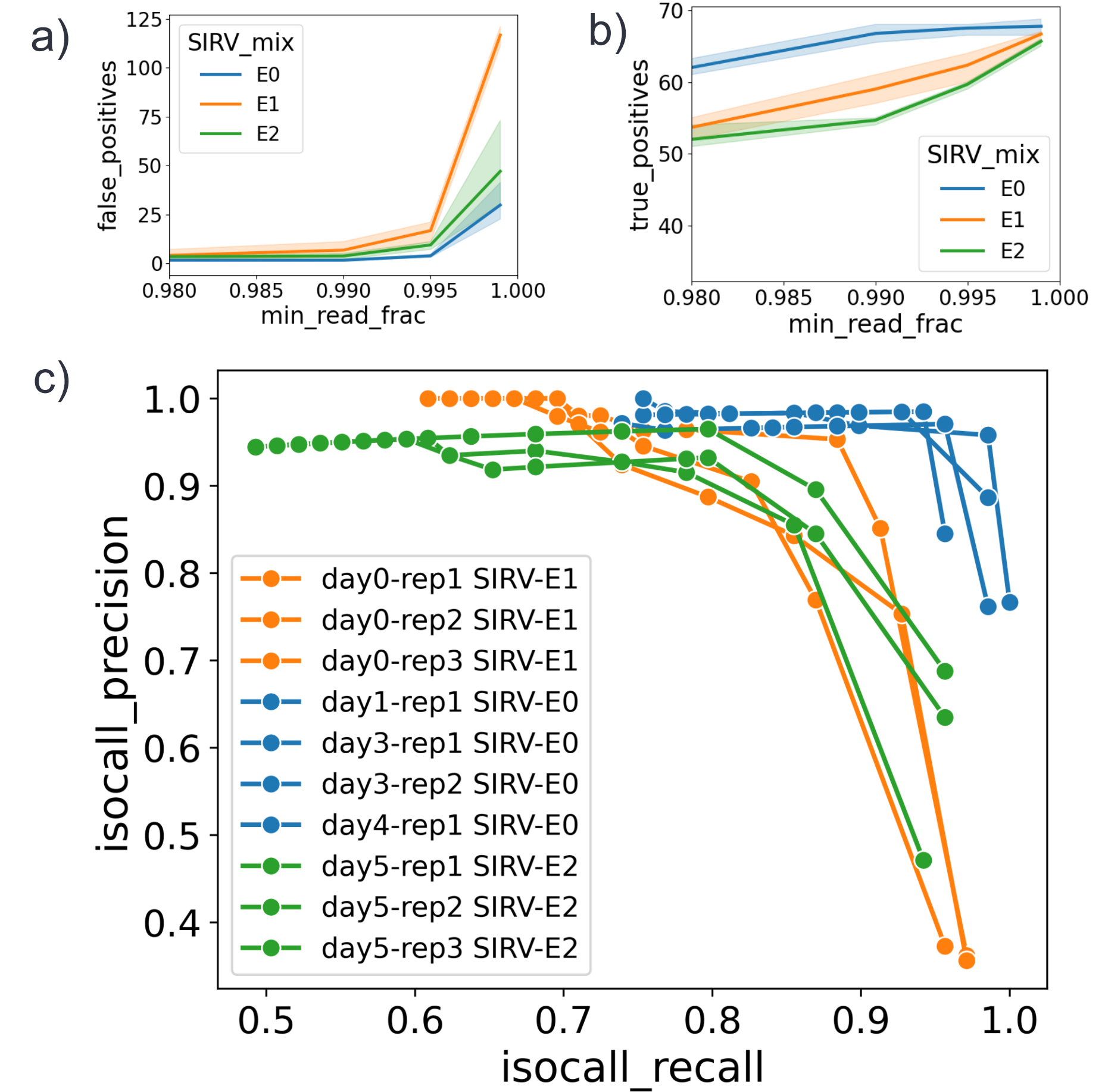


Figure 5. Sensitivity and specificity of isocall on SIRV spike-ins. a) b) increasing the minimum fraction of reads increases both true positives and false positives. c) precision and recall for each sample. Different SIRV mixes have different curves due to the relative abundances of the spike in transcripts.

Novel isoform on HEY2-AS

When looking at reads that support intron junctions, normalized to TPMs, we found novel isoforms with divergent abundance over time. One example is shown here.

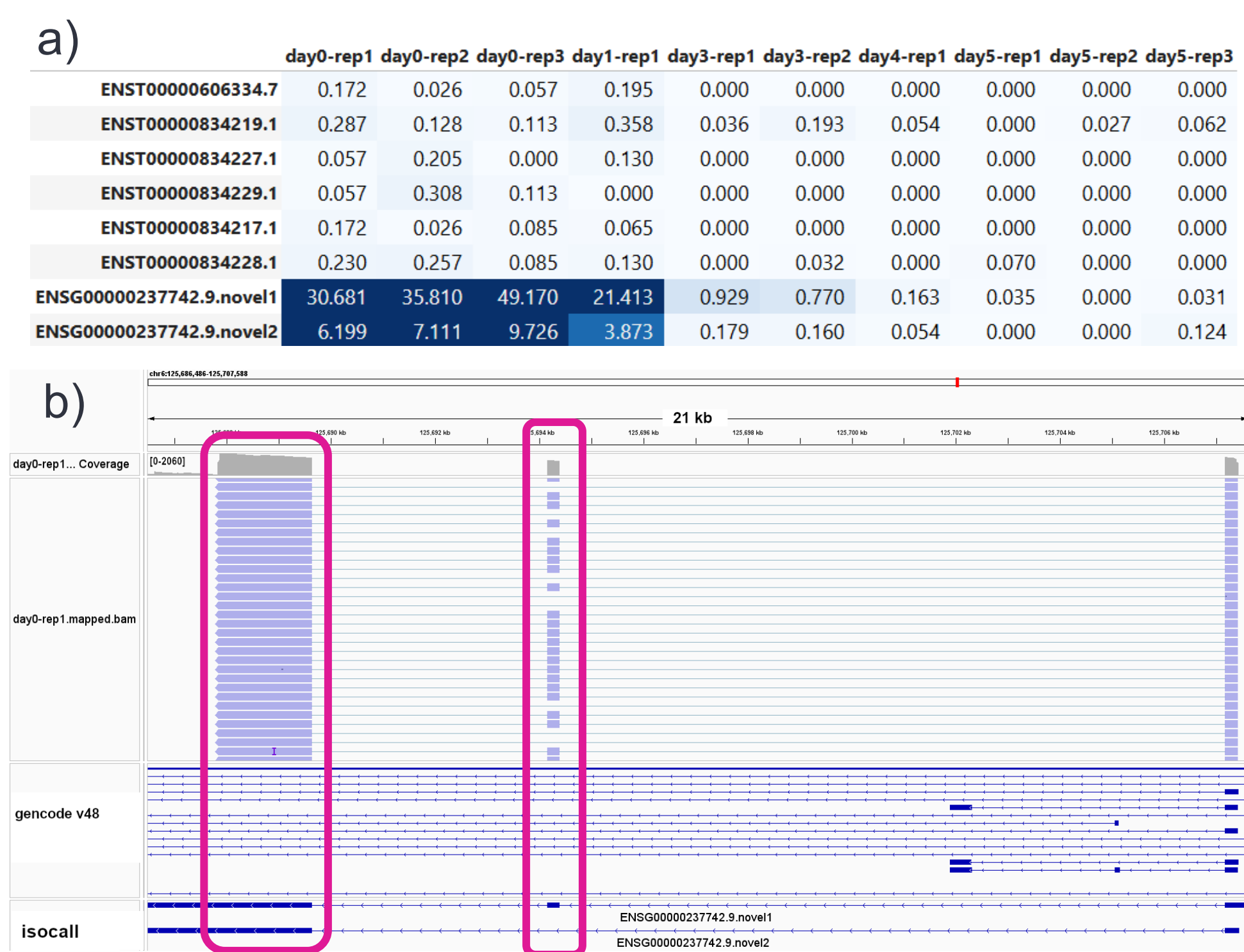


Figure 6. Novel isoform on HEY2-AS. a) CPMs for all known isoforms in gencode v48 on HEY2-AS and two most abundant novel isoforms. b) IGV screenshot of the novel isoforms and supporting reads. Highly expressed exons are not present in the reference annotation transcripts; many reads support the novel intron chains.

Conclusions

- Isocall enables population-scale joint isoform calling by analyzing hundreds of samples.
- Isocall allows for the discovery of novel isoforms.
- RNA workflows streamline analysis

References

1. A systematic benchmark of high-accuracy PacBio long-read RNA sequencing for transcript-level quantification

