

>12% of protein coding genes have >4-fold mRNA abundance difference depending on bioinformatic processing pipeline

Sonali Arora¹, Siobhan S. Pattwell¹, Eric C. Holland¹, Hamid Bolouri¹
Division of Human Biology, Fred Hutchinson Cancer Research Center, Seattle, WA 98109, USA



Abstract

Motivation: RNA-sequencing data is widely used to identify disease biomarkers and therapeutic targets. How reliable is it?
Approach: We compared data for 16730 genes from five “best-in-class” RNA-seq processing pipelines applied to 4,800 tumor and 1,890 normal human tissues from TCGA and GTEx.
Results: We show that for >12 % of protein-coding genes, RNA-seq expression estimates by different pipelines differ by > 4-fold in at least 10% of samples *using the same sequencing reads*.
Conclusions: A total of 2071 genes are discordantly quantified by current “best-in-class” RNA-seq processing pipelines.

Overview & Motivation

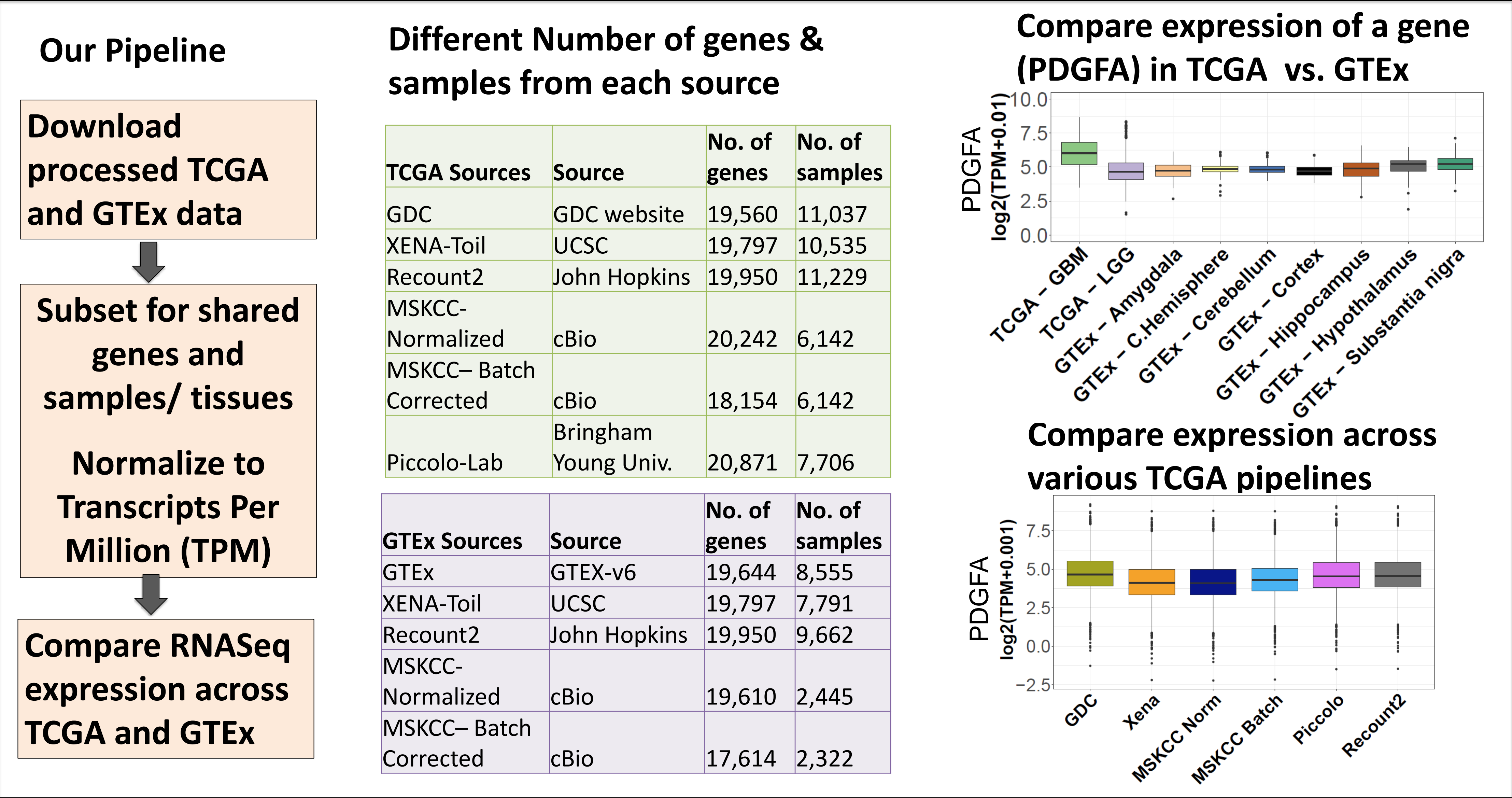


Fig 2: Expression estimates within individual pipelines are consistent, but differ across pipelines.

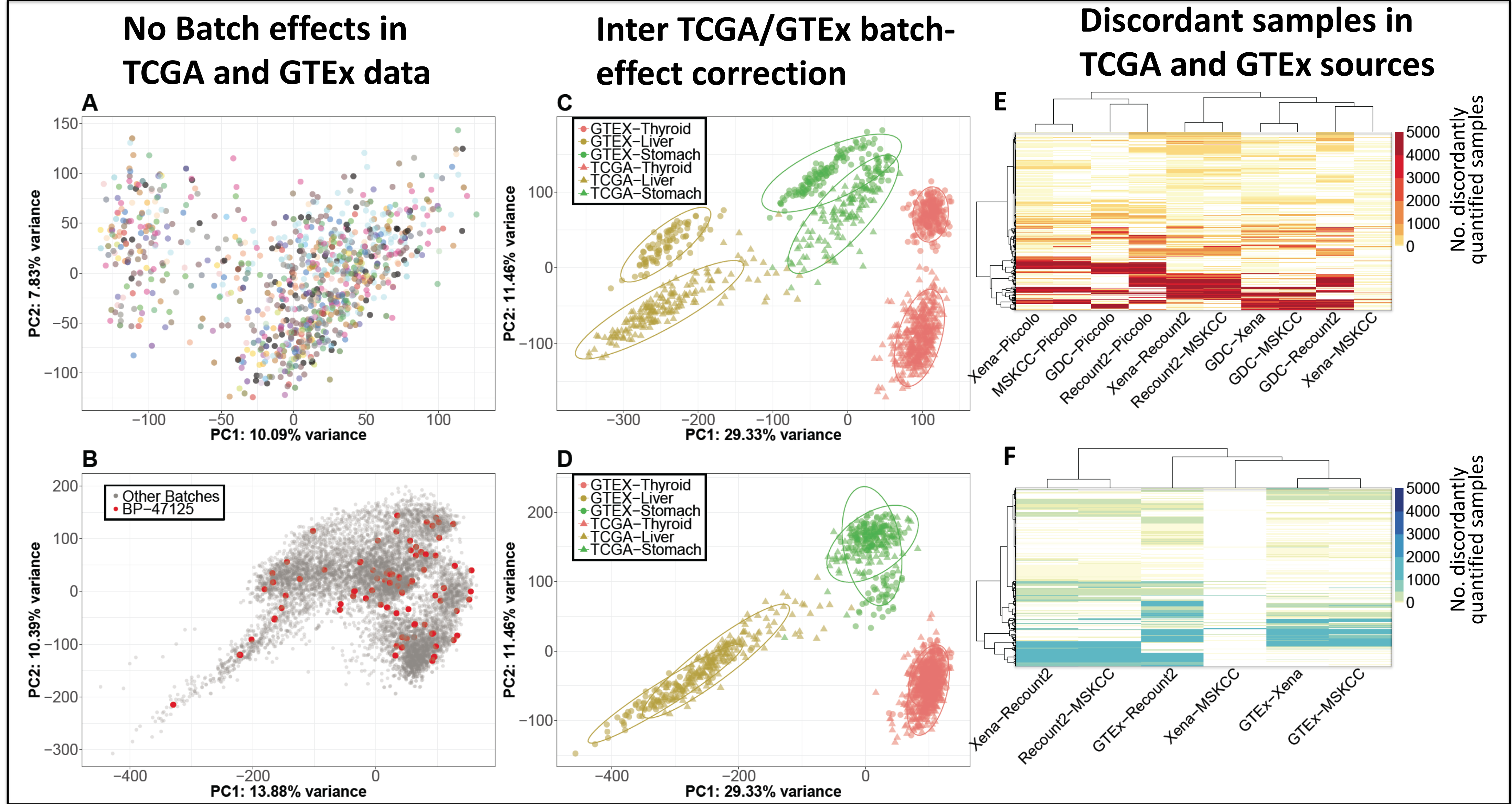


Fig 1: Normalization removes unwanted variability across Pipelines.

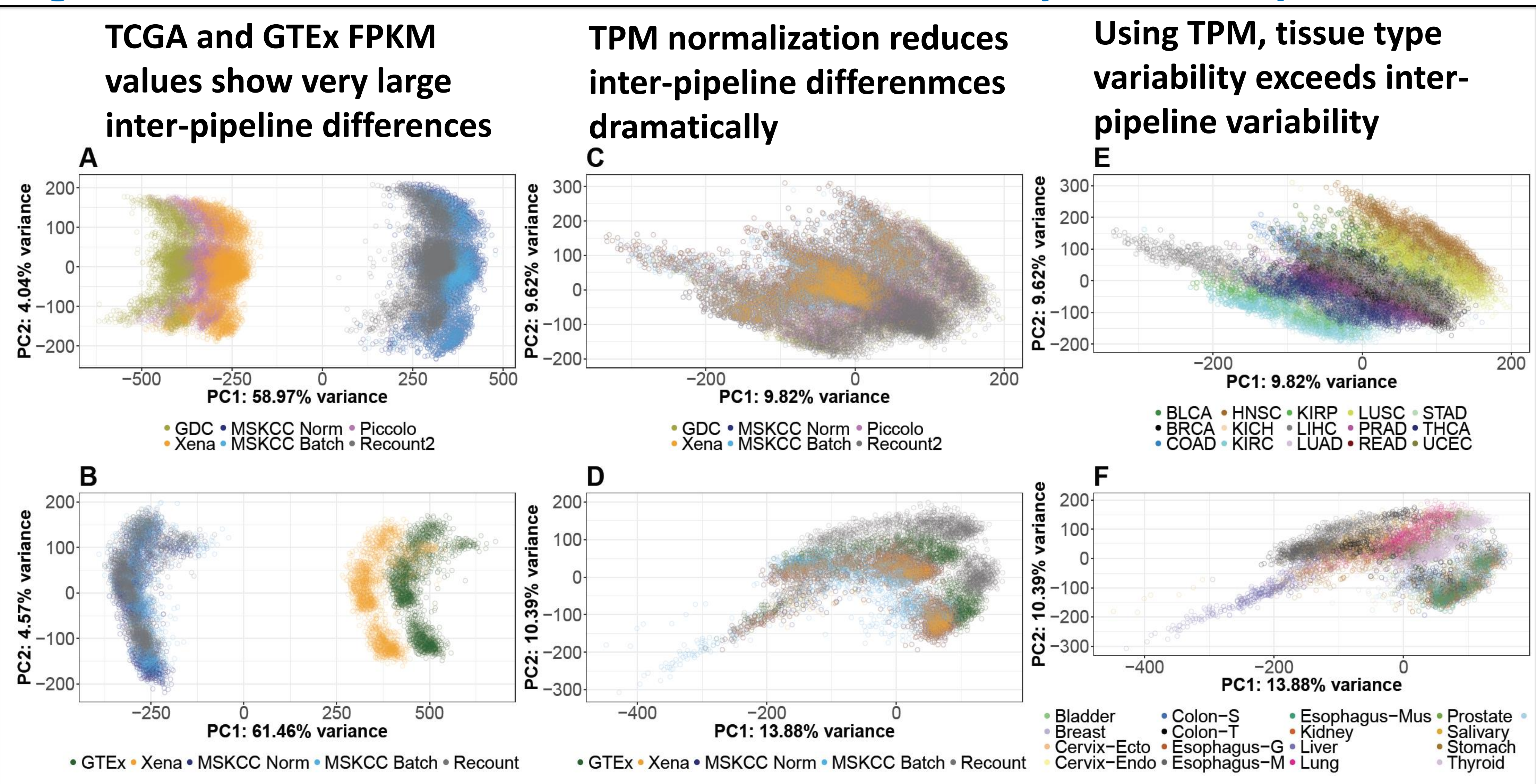
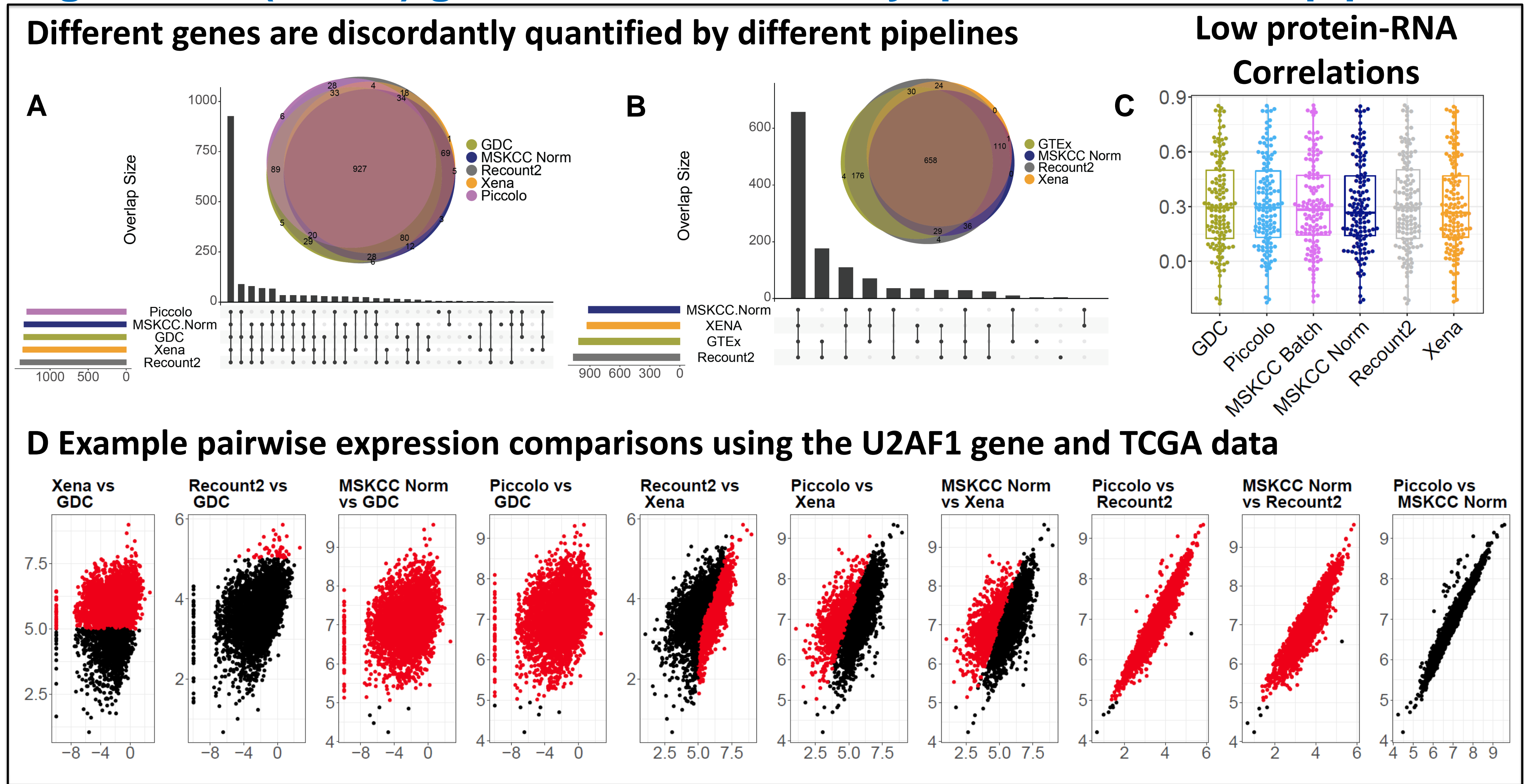


Fig 3: 2071 (12.4%) genes are discordantly quantified across pipelines.



References

1. Grossman, Robert L., Heath, Allison Pet al. (2016) Toward a Shared Vision for Cancer Genomic Data. New England Journal of Medicine
2. Vivian J, Rao AA, Nothaft FA, et al. (2017) Toil enables reproducible, open source, big biomedical data analyses. Nature biotechnology.
3. Collado-Torres L, Nellore A, et al (2017) Reproducible RNA-seq analysis using recount2. Nature biotechnology.
4. Rahman M, et al. (2015) Alternative preprocessing of RNA-Sequencing data in TCGA leads to improved analysis results. Bioinformatics.
5. Q. Wang, J Armenia, et al. (2018) Unifying cancer and normal RNA sequencing data from different sources. Scientific Data
6. The GTEx Consortium. The Genotype-Tissue Expression (GTEx) project. (2013) Nature genetics.

Example Impacted Cancer genes: ABCB6, ACY1, AKR1C1, AMACR, ANXA2, ARID4B, ARL2BP, ASPSCR1, ATF7IP, BMI1, BOP1, BRAF, BRD2, CD163L1, CEBPA, CNPY2, CSNK1A1, CSNK2B, DDR1, DLEC1, EIF4E, FAM168A, FCBGP, FHIT, GAGE1, GINS2, GNG11, GPR27, GSTM2, H3F3A, H3F3B, HIST1H1E, HIST1H3B, HOXD9, HRH2, KLRK1, KRT6B, KRT7, KYNU, LGALS7, LHFPL5, MIA, MORF4L1, MRPS18B, MUC6, MYCBP, P2RY8, PGAM1, PHLDB1, PKMYT1, PPIA, PSORS1C2, QKI, RHEB, RPL13, RPL14, RPL18, RPS3, RXRB, SLC19A1, SOD2, SPRR1A, STK19, STRADA, TATDN1, TNF, TPM3, TRIM27, TUBA1C, TXNDC17, UGT2B17, VDAC2, VPBEB1, WDR46, WFDC1, ZNF593, ZNF668