

Linking genomic variation to gene expression using pangenome graphs

Gözde YILDIZ

Email: goezde.yildiz@agrar.uni-giessen.de

Why a pangenome graph?

Sequence Variants

SNV (Single Nucleotide Variant)

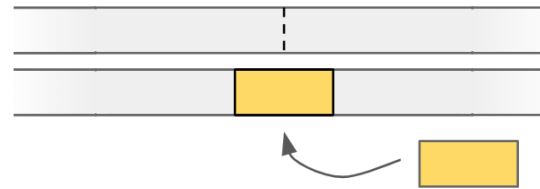


INDEL (Insertion or Deletion)

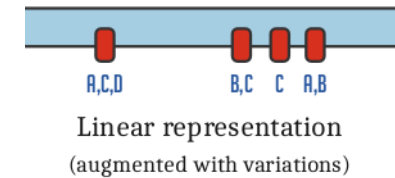
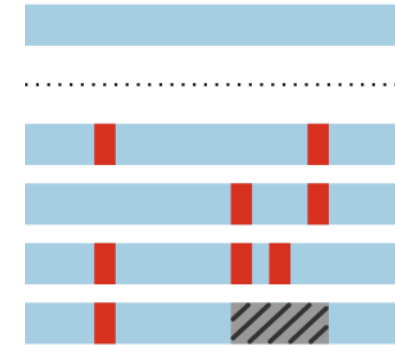


Structural Variants

Insertion



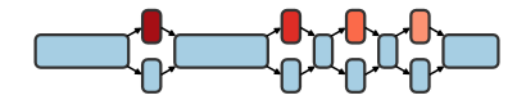
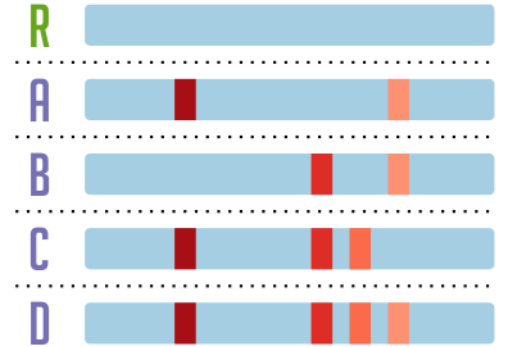
Inversion



Reference allele

Alternative allele

Unmapped segment



Graphical (compressed) representation

Shared segment

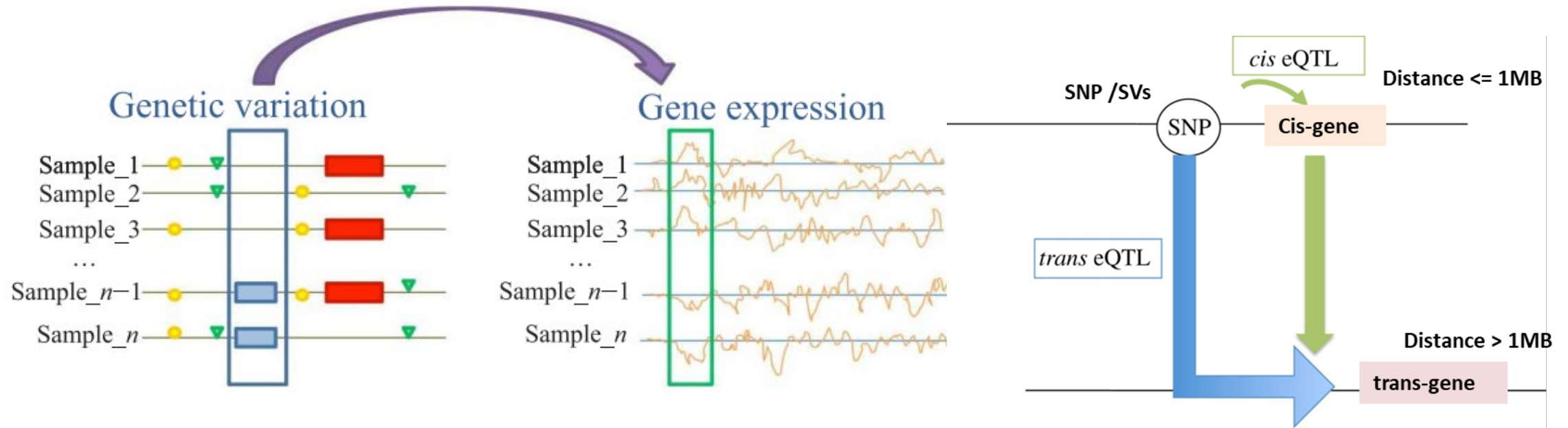
Variant

1. Encode **known variants** in a single data structure
2. **Avoid linear reference bias** during read mapping
3. **Improve** genotyping of SVs/SNPs

<https://pangenome.github.io/>

[https://www.melbournebioinformatics.org.au/tutorials/tutorials/longread sv calling/longread sv calling/](https://www.melbournebioinformatics.org.au/tutorials/tutorials/longread%20sv%20calling/longread%20sv%20calling/)

Expression quantitative trait locus (eQTL)

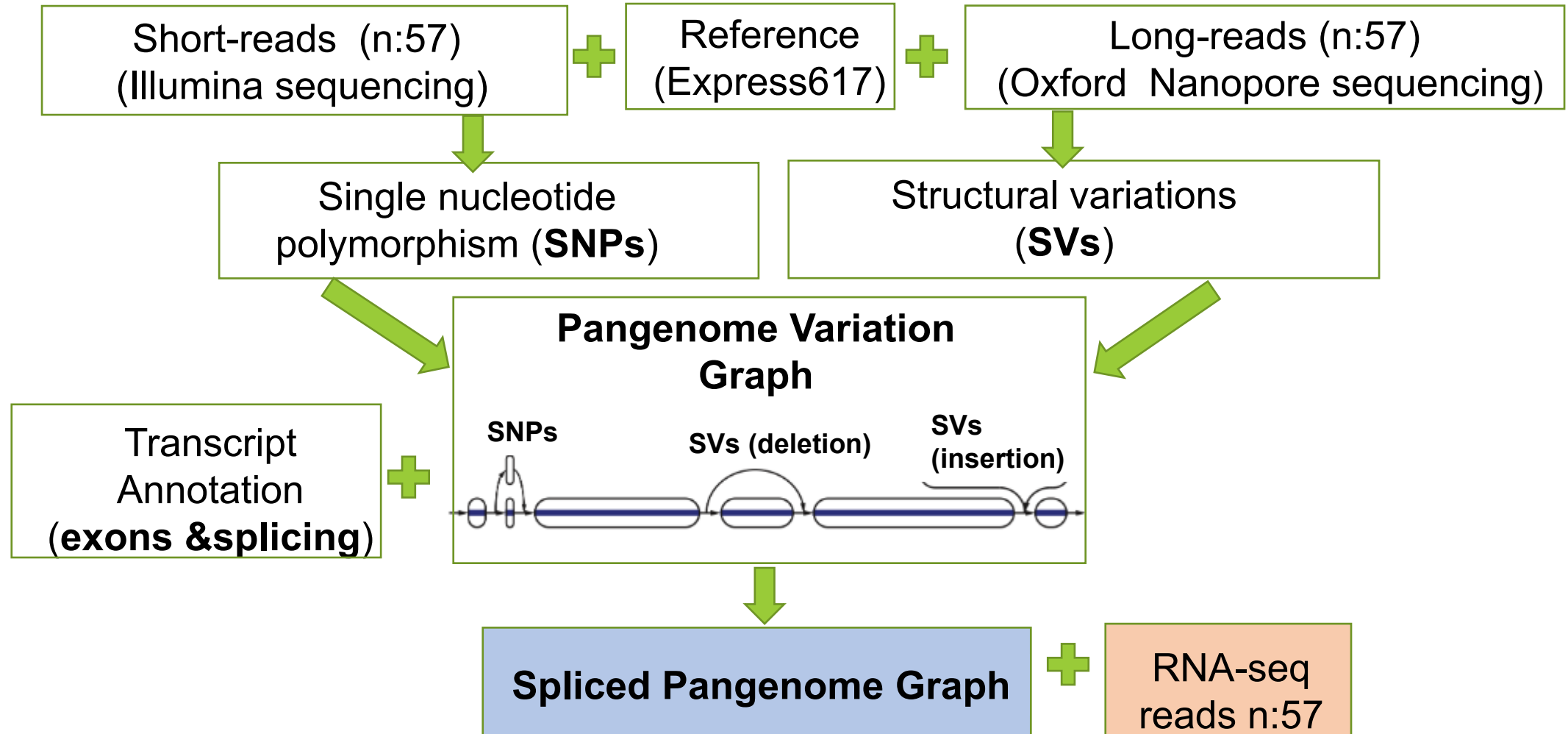


eQTL discovery from **linear-based** and **pangenome graph-based** approaches to assess the impact of SVs and SNPs on gene expression

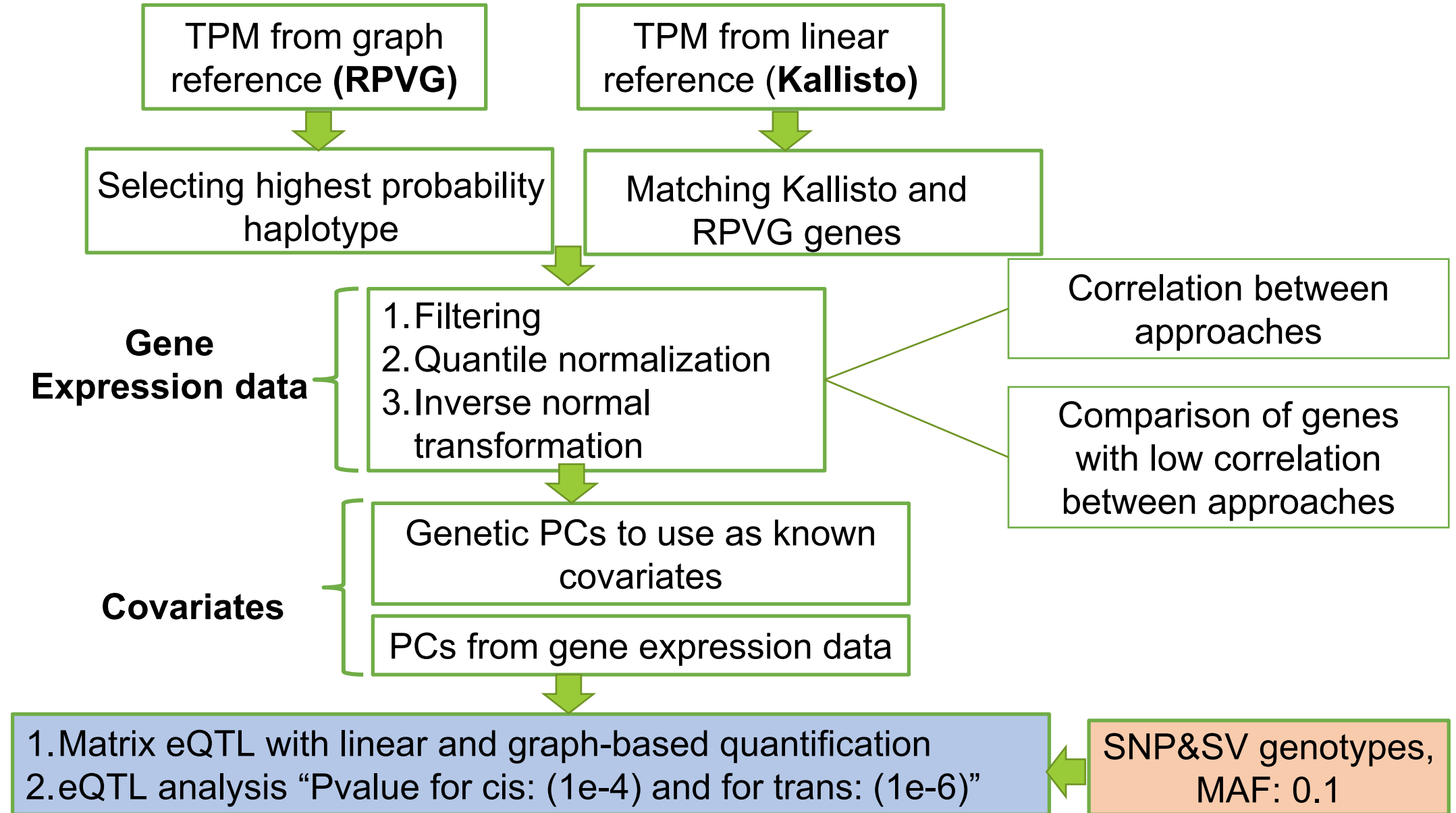
eQTL – finds genomic variations **statistically** associated with biological traits

[doi:10.1109/TST.2014.6961031](https://doi.org/10.1109/TST.2014.6961031)

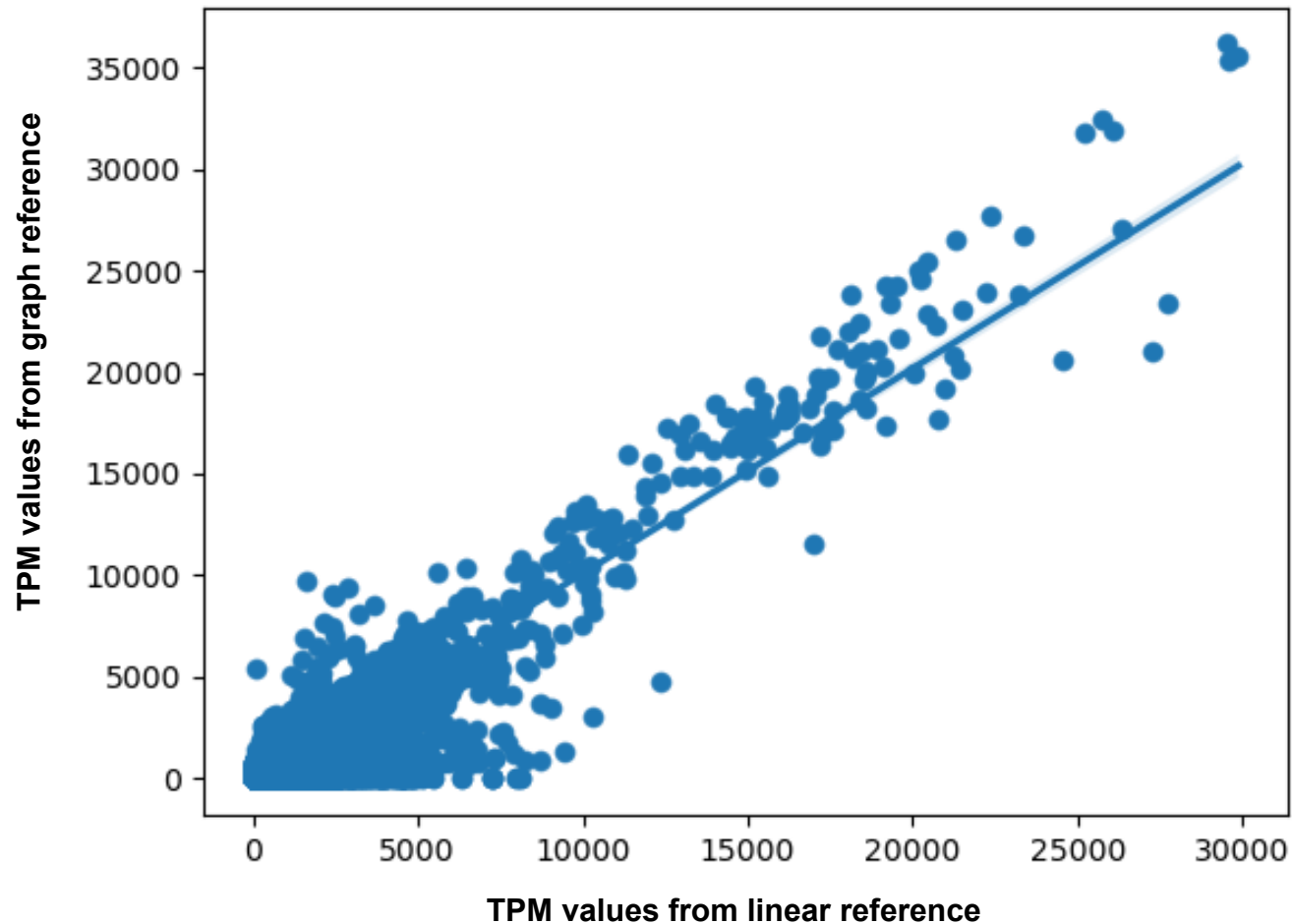
Building a variation graph for expression quantification



Designing eQTL analysis



TPM correlation between linear and pangenome graph references



Comparison of gene expression (TPM) from graph (rpvg, highest probability haplotype) and linear (Kallisto) references

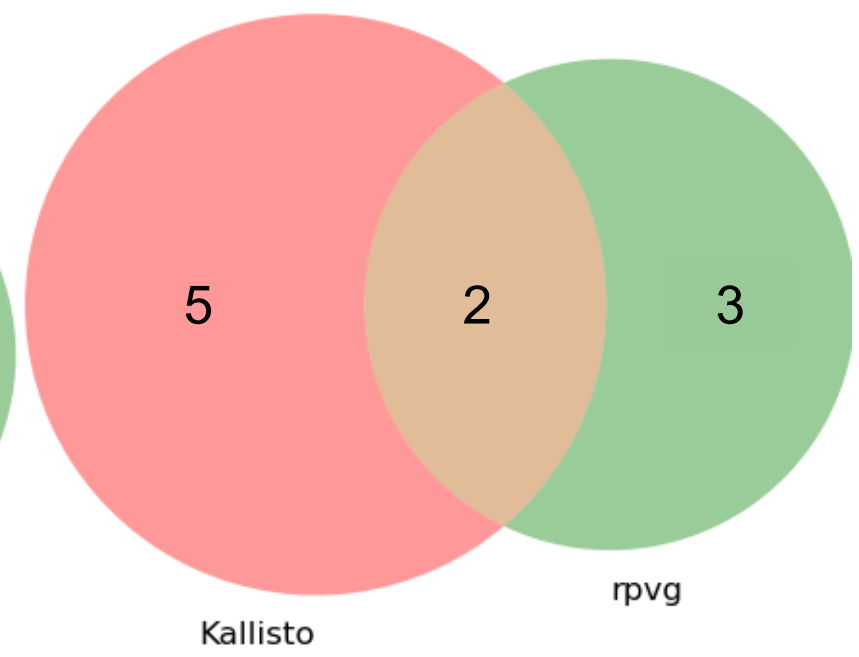
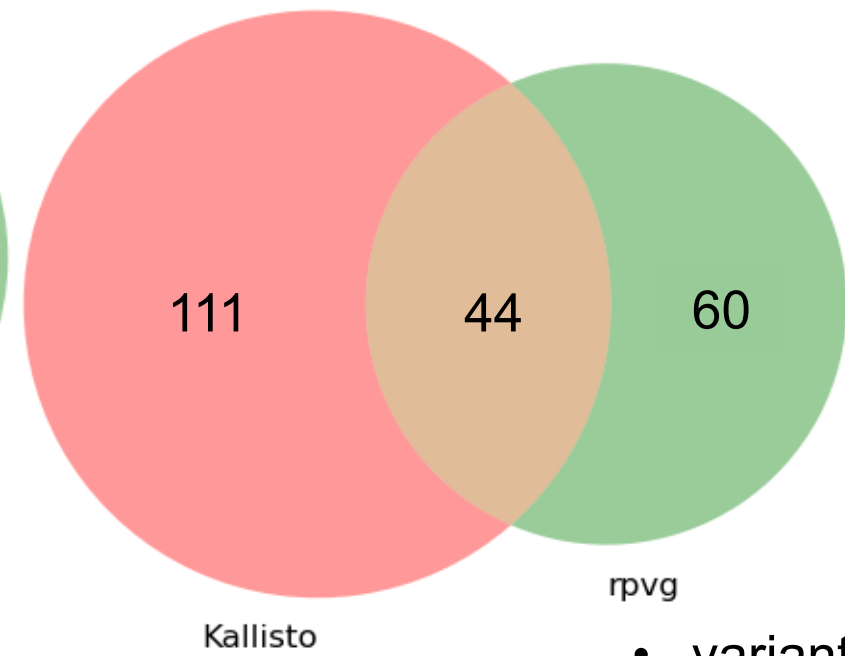
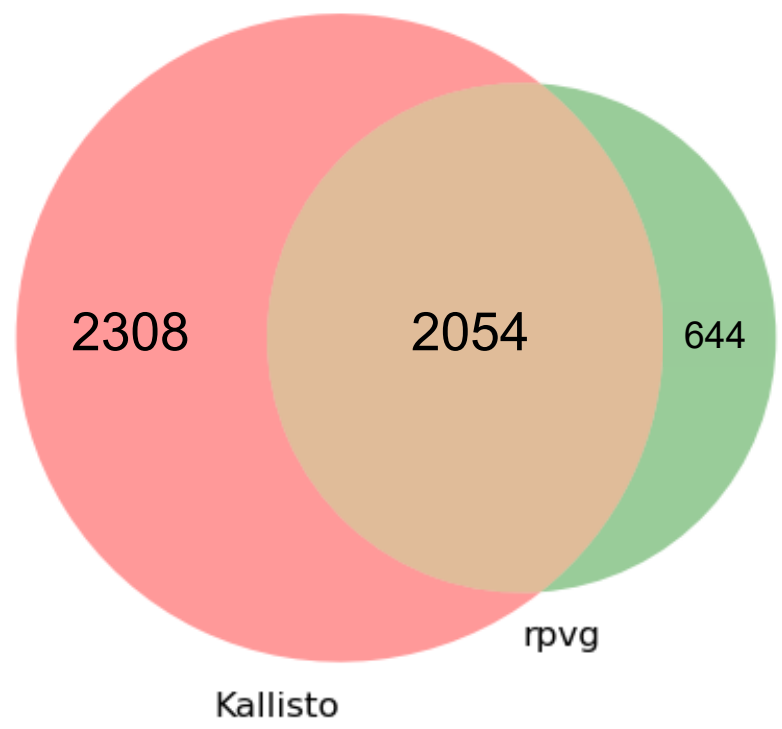
- **Pearson's r** : 0.938749954810782
- **Spearman's rho**: 0.9123681689144214
- **Kendall's tau**: 0.8091639040960162

Comparison of genes with eQTL identified by two approaches

Genes with SNP eQTL only

Genes with SV eQTL only

Genes with SNP & SV eQTL

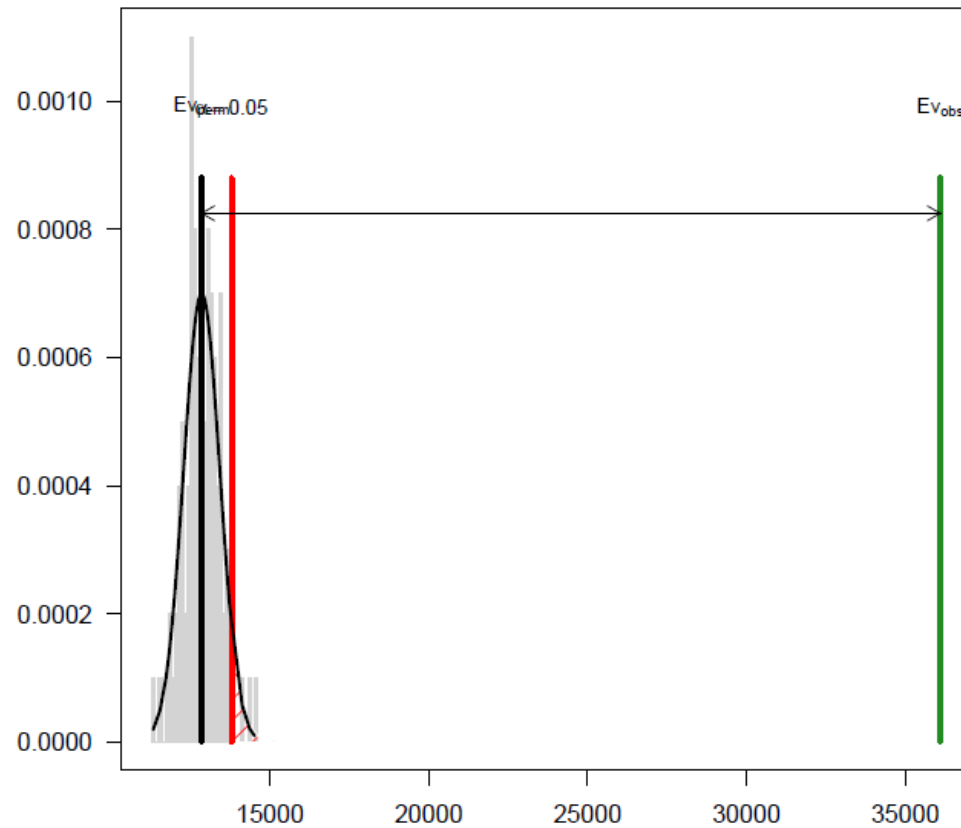


- variants with the lowest pvalues :
(4e-23) <= pvalues <= (1e-11) for Kallisto
(9e-33) <= pvalues <= (5e-12) for rpvg
- FDR threshold : **1e-6**

Genes with eQTL found by linear approach only are overrepresented in SNPs & SVs

SNPs

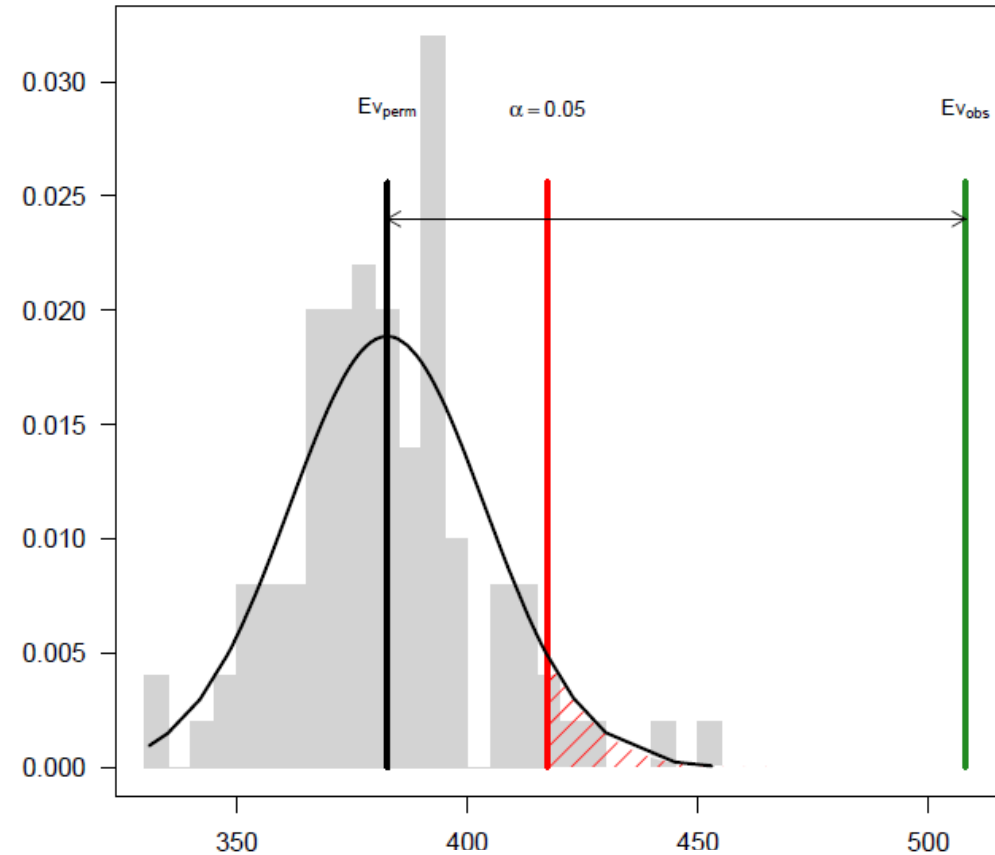
p-value: 0.0099
Z-score: 40.775
n perm: 100
randomization: resampleRegions



SNP number in genes

SVs

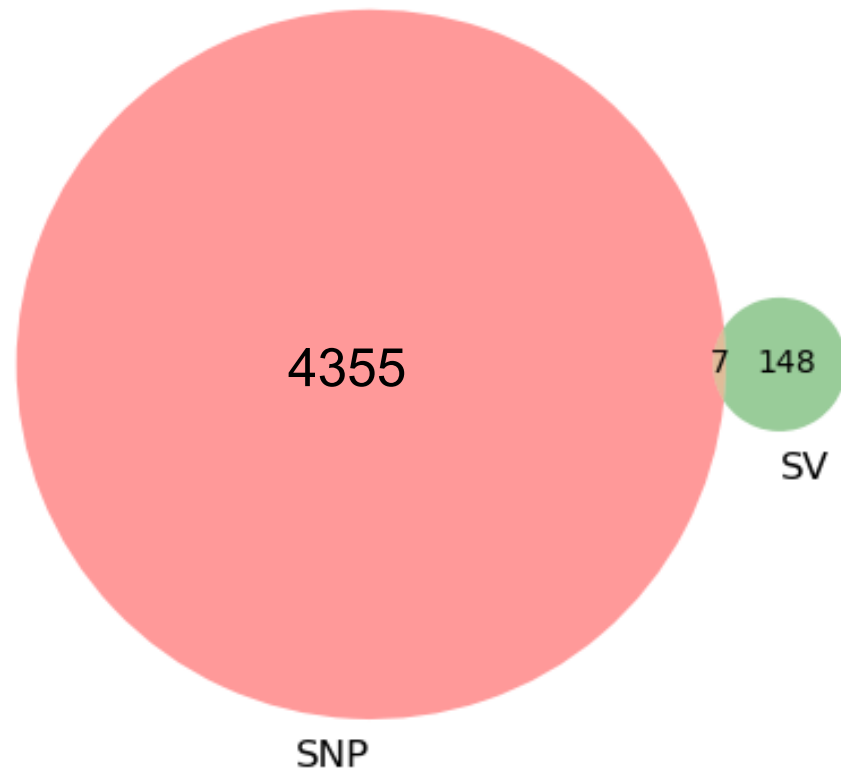
p-value: 0.0099
Z-score: 5.931
n perm: 100
randomization: resampleRegions



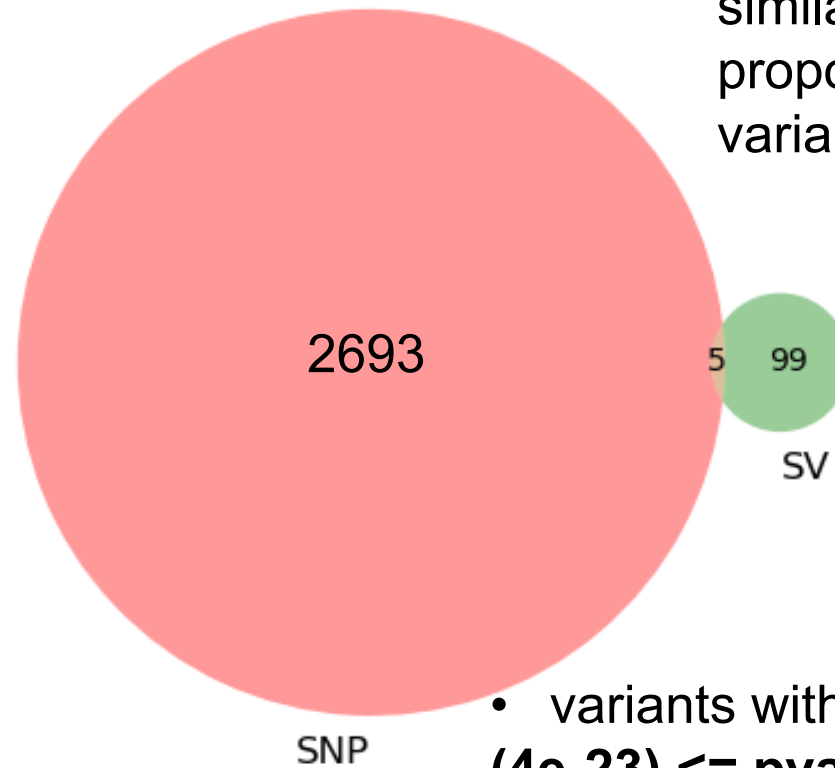
SV number in genes

Number of genes with SNP and SV eQTL

Linear based-approach



Graph based-approach



Proportion of eQTL SVs is similar to the overall proportion of SVs among all variants (**~3.5%**)

- variants with the lowest pvalues : **(4e-23) <= pvalues <= (1e-11)** for Kallisto
(9e-33) <= pvalues <= (5e-12) for rpvg
- FDR threshold : **1e-6**

Limitations and future work

- Limited samples size (57) allowing for detection of large effects only
- To expand dataset, SV genotyping from short reads can be performed (initial benchmarks suggest that ~50% of variants can be confidently genotyped)
- Functional analysis of eQTL genes under way

Acknowledgments

Agrobiinformatics group (JLU)

Dr. Agnieszka Golicz¹ (PI)

Dr. Silvia Zanini²

Jose Antonio Montero Tena³

Kevin Rockenbach⁴

Kübra Arslan⁵

Venkataramana Kopalli⁶

Rishi Srivastava⁷



goezde.yildiz@agrار.uni-giessen.de

agnieszka.golicz@agrار.uni-giessen.de



Plant Breeding Department (JLU)

Prof. Dr. Rod Snowdon⁸