**a)**

**Manuscript title**

> Authors

∨ Abstract

∨ Introduction

∨ Results

**PhenoPLIER: an integration framework based on gene co-expression patterns**
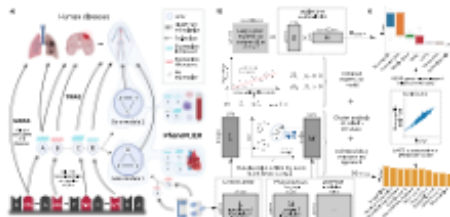
Figure 1: Schematic of the PhenoPLIER framework. a) High-level schematic of PhenoPLIER (a gene module-based method) in the context of TWAS (single-gene) and GWAS (genetic variants). PhenoPLIER integrates groups of genes co-expressed in specific cell types (gene modules) with gene-trait and gene-drug associations.

∨ Discussion

∨ Methods and materials

$$y = t_l \gamma_l + e_l,$$ (1)

$$\hat{z}_l = \sum_{a \in model_l} w_a^l \frac{\hat{\sigma}_a}{\hat{\sigma}_l} \frac{\hat{\beta}_a}{se(\hat{\beta}_a)},$$ (2)

**Code and data availability**

**b) Section-specific prompt generator:**

Revise the following paragraph from the abstract of an academic paper (with the title 'Manuscript title' and keywords 'keyword1, keyword2, ...') so the research problem/question is clear, the solution proposed is clear, the text grammar is correct, spelling errors are fixed, and the text is in active voice and has a clear sentence structure.

Revise the following paragraph from the Introduction... so most of the citations to other academic papers are kept, the text minimizes the use of jargon, ...

Revise the following paragraph from the Results... so most references to figures and tables are kept, the details are enough to clearly explain the outcomes, sentences are concise and to the point, ...

Revise the following paragraph from the Discussion... so most of the citations to other academic papers are kept, the text minimizes the use of jargon, ...

Revise the following paragraph from the Methods... so most of the citations to other academic papers are kept, most of the technical details are kept, most references to equations (such as "Equation (@id)") are kept, all equations definitions (such as "$$ ... $$ {#id}") are included, the most important symbols in equations are defined, ...

...