
Small-cohort GWAS discovery with AI over massive functional genomics knowledge graph

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Genome-Wide Association Studies (GWAS) links genetic markers with diseases
2 and is the cornerstone for the development of effective therapeutics. However, for a
3 long tail of many uncommon diseases, the small GWAS sample sizes limit detection
4 power and hamper development of effective treatments. The recent substantial
5 growth in the size of functional genomics data presents a fresh opportunity to
6 tackle these challenges. Here, we introduce KGWAS, a novel geometric deep
7 learning method that leverages a knowledge graph to integrate massive functional
8 information about variants, genes, gene programs, and their interactions, assessing
9 variant-disease associations. Unlike conventional GWAS, which treats variants
10 independently, our approach recognizes that variants influence disease through
11 complex cellular networks. Our realistic simulations show that KGWAS is well-
12 calibrated and powerful in identifying disease variants. We applied KGWAS to
13 21 independent UK Biobank diseases/traits from small subsampled cohorts ($N=1$ -
14 10K), and KGWAS produced significantly more independent associations that
15 were replicable in the full cohort (average $N=374K$), 22.0%-89.9% higher than
16 state-of-the-art baselines. Next, we applied KGWAS to 554 less common UK
17 Biobank diseases ($N_{\text{case}} < 5K$) and identified 183 novel loci, 46.9% higher than the
18 original GWAS, including rs2155219 associated with ulcerative colitis potentially
19 via regulating *LRRC32* expression in CD4+ regulatory T cells, and rs73127651
20 associated with myasthenia gravis potentially via regulating *PPHLN1* expression in
21 brain cell types. Overall, KGWAS is a flexible and powerful AI model to integrate
22 the growing functional genomics data to discover novel variants for small cohort
23 diseases.