

Deep Learning Based Phenotyping of Medical Images Improves Power for Gene Discovery of Knee Osteoarthritis

Eugenia Lin, Brianna I Flynn, Emily M Javan¹, Zoe Trutner, Karl Marc Koenig, Kenoma O Anighoro, Alaukik Gupta², Prakash Jayakumar, Vagheesh Narasimhan

¹Department of Integrative Biology, ²Department of Biomedical Engineering

INTRODUCTION: Recent genetic studies have successfully applied deep-learning methods to generate image derived phenotypes (IDPs) of various disease states and linked them with genome-wide significant loci. While some recent studies on musculoskeletal disease employ these novel phenotyping approaches, to our knowledge, no studies have investigated how generating quantitative IDPs that underlie binary disease status could be used to improve power for gene discovery at the population scale. For radiographically diagnosed diseases, such as knee osteoarthritis (OA), computer vision approaches for automated phenotyping based on training data from clinicians offer the potential to ascertain both case status and disease severity at scale. In this study, we thus sought to analyze the genetic basis of knee OA associated phenotypes within the United Kingdom Biobank (UKB).

METHODS:

We first trained a deep learning binary model (DL-binary) to identify knee OA cases from the level of individual clinicians and deployed this model on anterior-posterior (AP) knee radiographs at biobank scale of 28,725 individuals in the UKB, all of whom have genome sequence data. We then compared our radiographically obtained results to the ICD-10 record. Second, we trained a separate image segmentation algorithm to obtain a quantitative measurement highly correlated with knee OA severity, minimum joint space width (mJSW), to examine differences in power between genome wide association studies (GWAS) carried out using quantitative approaches versus a case-control design. Finally, we generate a polygenic risk score (PRS) for each phenotype to evaluate if improvements in statistical power to find novel loci translate to better prediction of ICD-10 record knee OA (M17) in a hold-out dataset of over 300,000 individuals.

RESULTS:

Using our model, we identified 2,603 (240%) more cases than currently diagnosed with M17 (knee OA) in the ICD-10 record. Individuals diagnosed as cases by the DL-binary model had higher rates of self-reported knee pain, for longer durations and with increased severity compared to control individuals. Despite the DL-binary and mJSW phenotypes being highly genetically correlated (92%), the heritability of mJSW ($h^2_g=0.24$) was an order of magnitude greater than that of using ICD-10 codes to identify knee OA ($h^2_g=0.02$) or DL-binary ($h^2_g=0.04$) phenotypes. In a GWAS run on the mJSW phenotype, we identified 18 genome-wide significant loci, as opposed to 1 and 6 at the same sample size using DL-binary and or ICD-10 coding, respectively. The improved power of the quantitative phenotype also translated to better polygenic risk score (PRS) prediction for knee OA diagnosis in a holdout dataset of 371,686 individuals.

DISCUSSION AND CONCLUSION:

In this study, we demonstrate a deep learning method to directly phenotype OA cases and controls (the DL-binary model), as well as joint space narrowing (mJSW-model), from AP view knee radiographs of the UKB. We compared this image derived phenotyping approach with case-control status already available in the ICD-10 record code on the same set of individuals, to determine if image-derived phenotyping approaches have an effect on statistical power in GWAS. We find that the case-control phenotyping using the DL-binary classification method enables us to raise the case count by greater than two fold. While computer vision approaches to extract and analyze radiograph derived phenotypes are not themselves novel, this work is among the first to use this approach on a disease for which diagnosis is primarily radiographic, to demonstrate that having a quantitative endophenotype that captures additional information about variation in disease severity improves power for genomic and epidemiological analysis.

