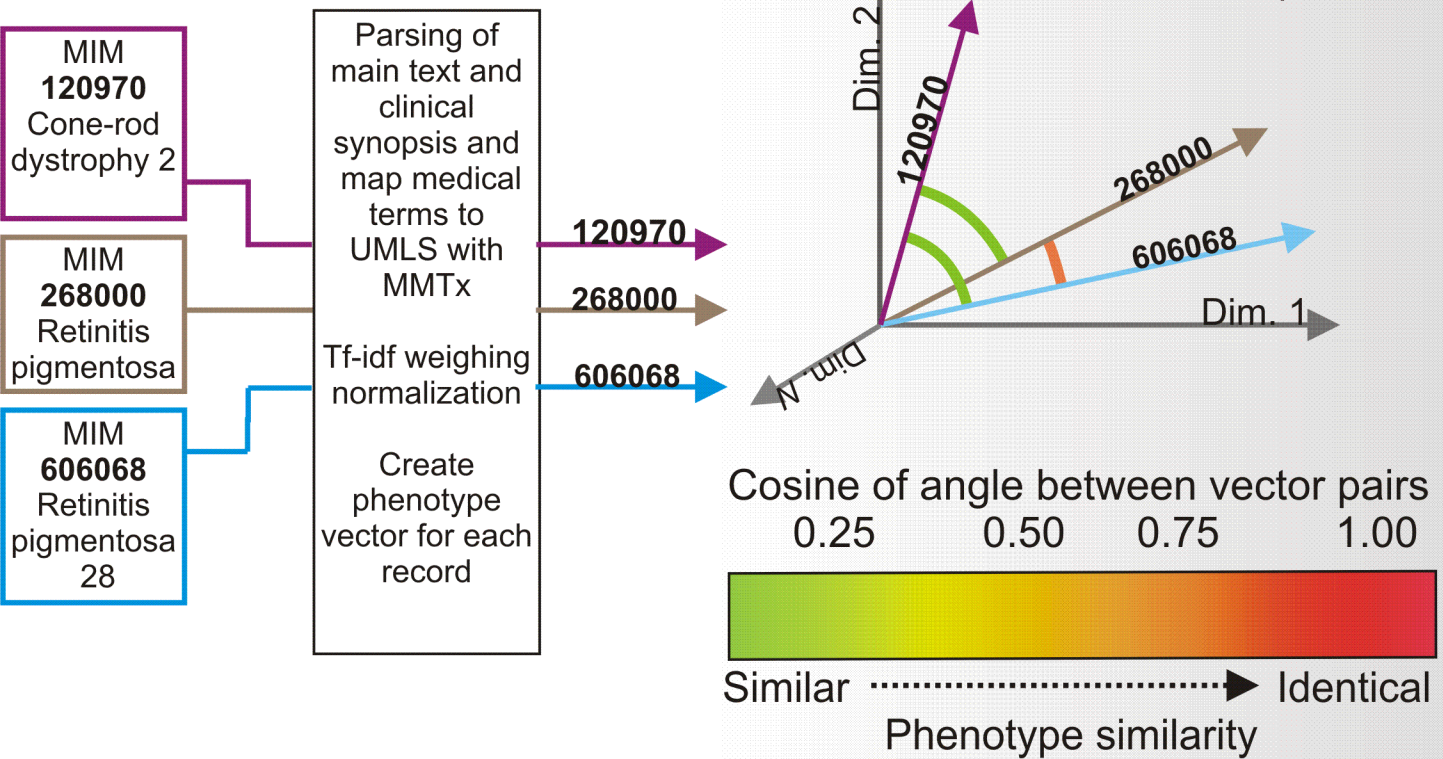
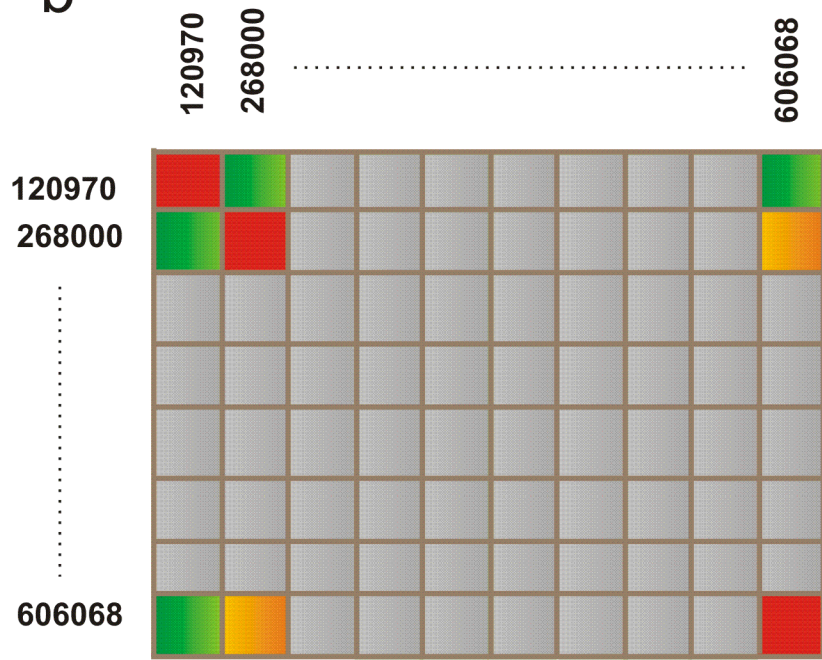


a**b**

Supplementary Figure 1 Measuring phenotype association scores between OMIM records. For every record we created a phenotype vector **a**). A phenotype vector consists of weighted medical terms present in the record, and represents the phenotype described in that particular record. For a single record a vector was made by parsing the clinical synopsis and main text of a record with MMTx. MMTx maps the words to medical terms in a subset of the UMLS vocabulary. Every term in the vector was weighted in relation to the record it was describing, and the weight normalizes the term scores in relation to the length of the records. Based on these medical terms a vector representing every OMIM record was then plotted in medical term space. The pairwise phenotypic overlap between records was quantified by the cosine of the angle between normalized vector pairs. Doing this for all combinations of vector pairs, enabled us to make a matrix of pairwise phenotypic similarity scores between all OMIM records **b**). For clarity we only show the procedure for three OMIM records (MIM 120970, MIM 268000, MIM 606068).