# Supplemental information

# Identification of risk genes for

# Alzheimer's disease by gene embedding

**Yashwanth Lagisetty, Thomas Bourquard, Ismael Al-Ramahi, Carl Grant Mangleburg, Samantha Mota, Shirin Soleimani, Joshua M. Shulman, Juan Botas, Kwanghyuk Lee, and Olivier Lichtarge**

1 **Supplementary Information for**

2 Identification of Risk Genes for Alzheimer's Disease by Gene

3 Embedding.

4

5 Yashwanth Lagisetty[1,2], Thomas Bourquard[2], Ismael Al-Ramahi[2,3,4], Carl Grant Mangleburg[2],

6 Samantha Mota[2], Shirin Soleimani[2], Joshua M. Shulman[2,3,4,5,6], Juan Botas[2,3,4], Kwanghyuk Lee[2],

7 Olivier Lichtarge* [2,4,7]

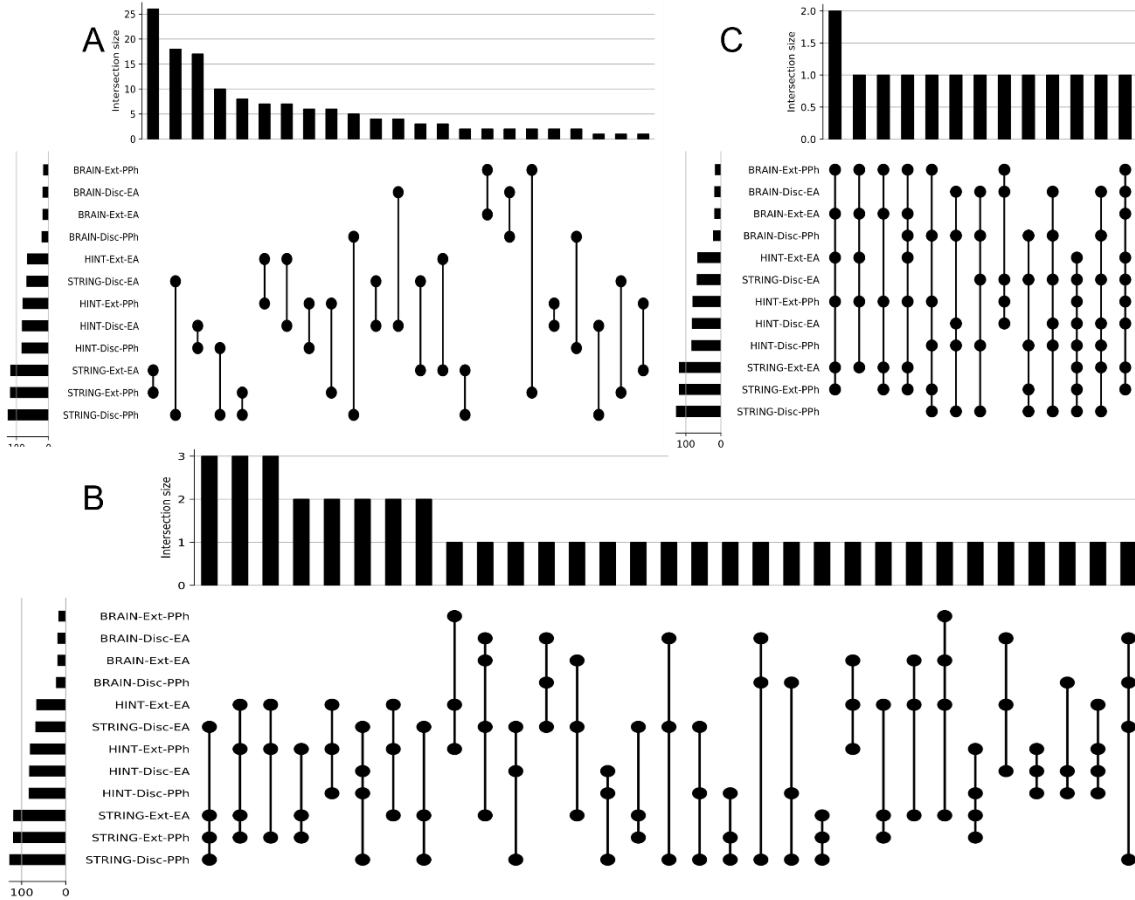8 * To whom correspondence should be addressed:

9 Olivier Lichtarge MD, PhD

10 **Email:** lichtarge@bcm.edu
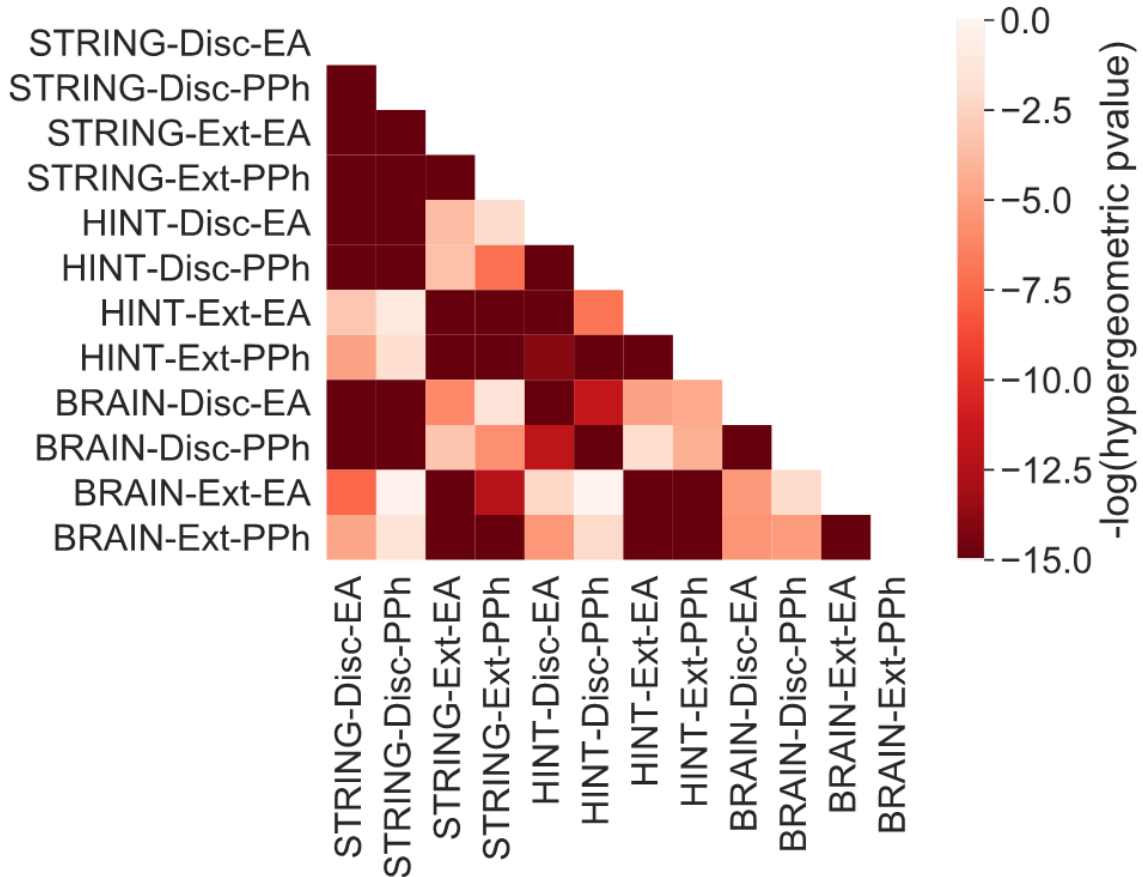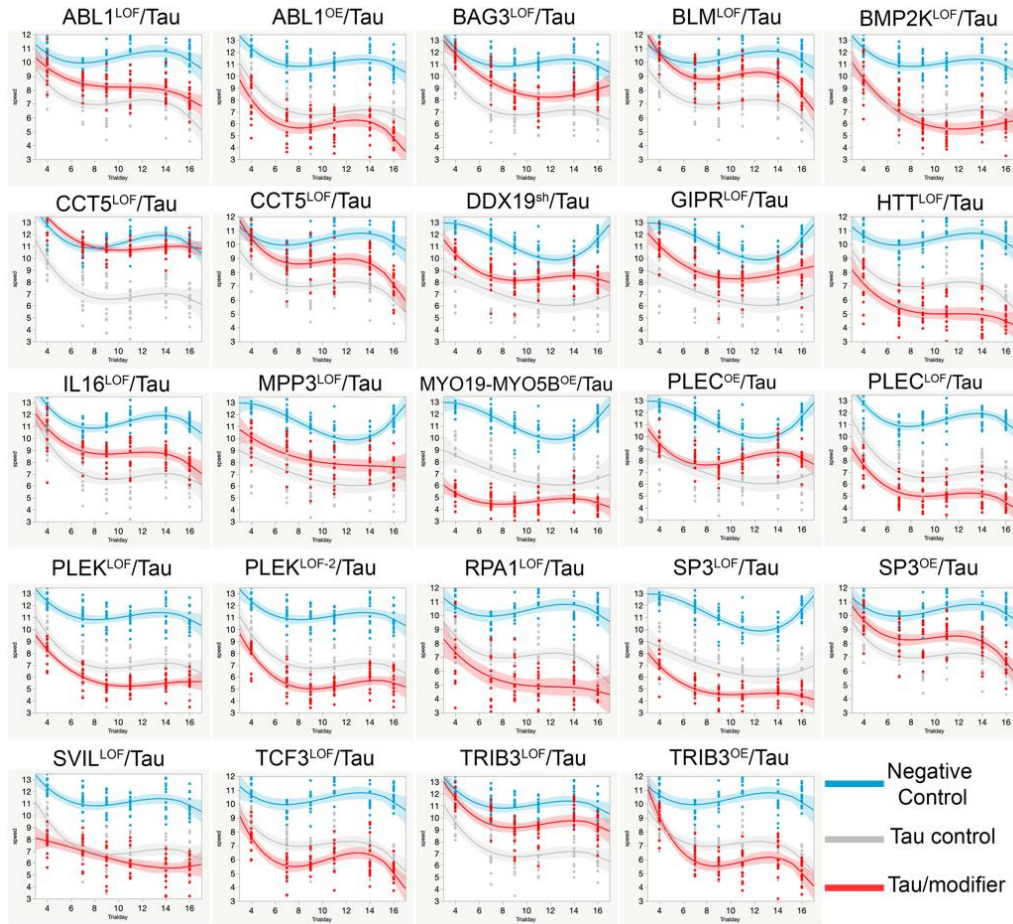
11

12

13 **Supplemental Figure Legends**



14

15 **Figure S1. UpSet plot of combinatorial intersections between all 12 GeneEMBED**

16 **experiments, Related to Figure 4.** (A) pairwise intersections. (B) Intersections between 3-4 sets.

17    (C) Intersections between 5+ sets. Set of 143 unique genes from all intersections are used for 'high-

18    confidence' gene set.

19



20

21    **Figure S2. GeneEMBED candidates are consistently identified across various cohorts,**

22    **networks, and VIS systems, Related to Figure 4.** One-tailed hypergeometric overlap tests were

23    done on every pairwise combination of cohort-network-VIS experiments. Among 66 independent

24    pairwise tests, only 11 did not demonstrate statistically significant hypergeometric p-values ($p <$

25    0.05, log(p) < -2.99).

26

**Figure S3. GeneEMBED candidates modulate tau-induced neuronal dysfunction, Related to Figure 4.** Regressions representing average speed as a function of age in control fruit flies (blue) or flies expressing human wild type _Tau_ either alone (grey) or together with the above indicate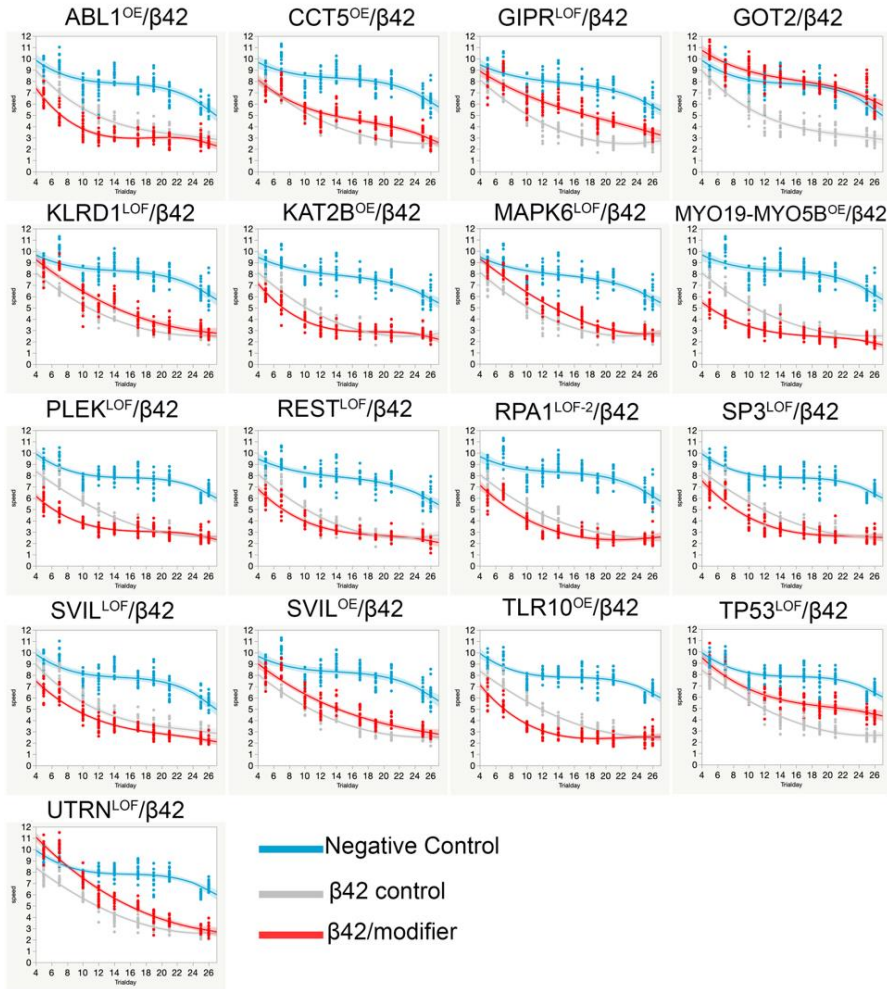d modifiers (red) on the corresponding Drosophila homolog (see supplementary table 12 for genotype details). Charts show third degree polynomials and confidence intervals. All differential effects were statistically significant ($p<0.01$) following ANOVA analysis on Linear mixed models regression with fitted splines

36

**Figure S4. GeneEMBED candidates modulate β amyloid-induced neuronal dysfunction, Related to Figure 4.** Regressions representing average speed as a function of age in control fruit flies (blue) or flies expressing human wild type _β amyloid_ either alone (grey) or together with th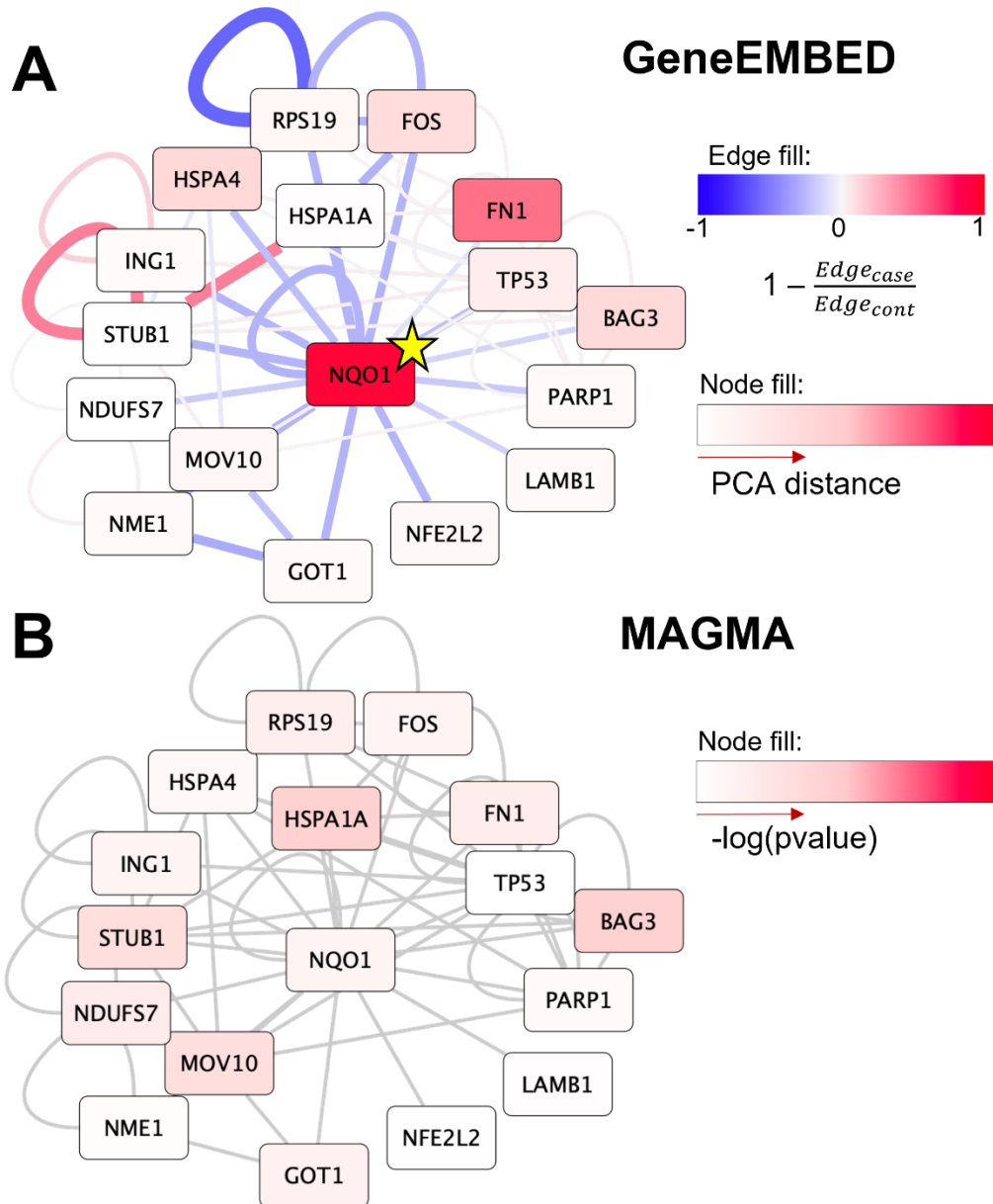e above indicated modifiers (red) on the corresponding Drosophila homolog (see supplementary table 12 for genotype details). Charts show third degree polynomials and confidence intervals. All differential effects were statistically significant (p<0.01) following ANOVA analysis on Linear mixed models regression with fitted splines

**A** GeneEMBED

**B** MAGMA

**Figure S5. Visual example of GeneEMBED's network informed gene discovery, Related to Figure 1.** (A) Network of *NQO1* from the Brain network. Edge color represents the zero-centered ratio of mutation edge weight in c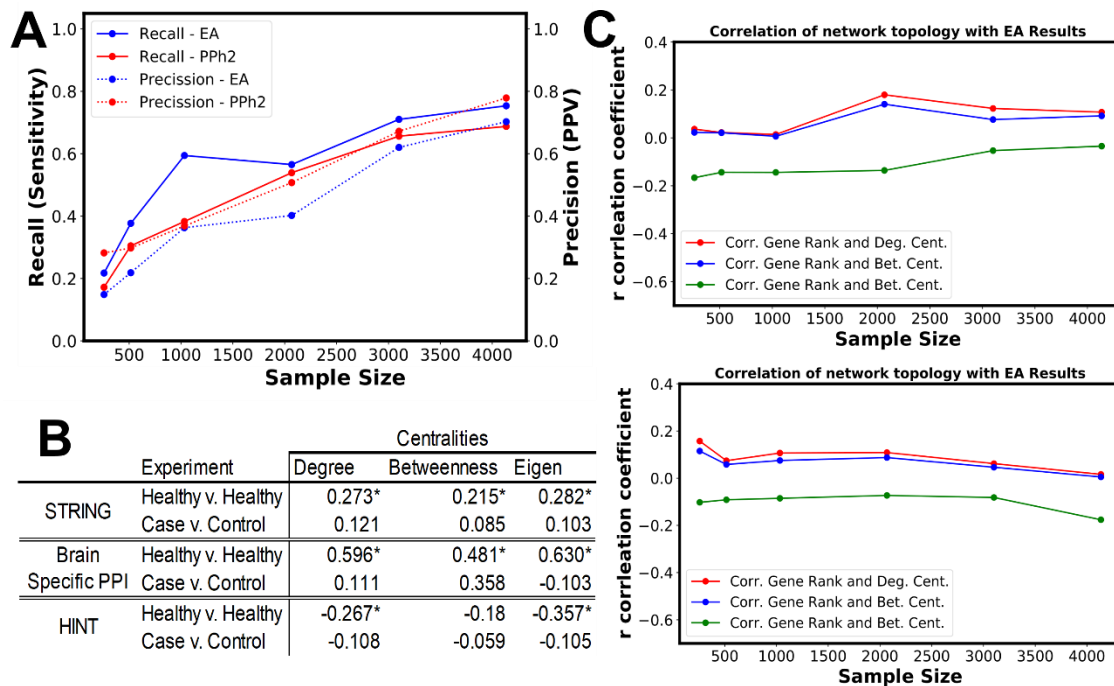ases versus controls. Edge width represents the magnitude of this ratio. Node fill is represented by PCA distance from GeneEMBED on the Discovery cohort using EA. The star on *NQO1* indicates that this gene was identified with FDR < 0.01 in GeneEMBED analysis. (B) shows the same network but with node fill corresponding to the -log(pvalue) from MAGMA analysis on the Discovery cohort. Subtle network differences allow GeneEMBED to identify *NQO1* when mutational data alone would not suffice.
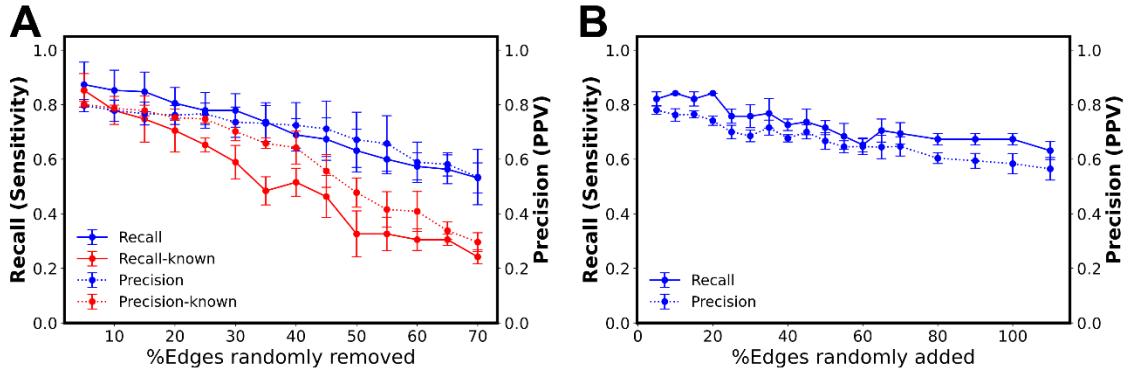
54



55

56 **Figure S6. GeneEMBED is robust to low sample sizes, Related to Figure 1.** (A) Plot of precision

57 and recall of GeneEMBED identified genes at decreased sample sizes relative to genes identified

58 using the full Discovery cohort. (B) Spearman rank-order correlation between genes identified using

59 the three brain networks applied to Healthy vs Healthy controls or case vs control experiment.

60 Asterisk indicates statistically significant (p<0.05) correlation. When disease relevant information is

61 removed from data, GeneEMBED relies on network topology to rank genes. (C) Spearman rank-

62 order correlation between candidates identified at low cohort sizes.

63

64

65

**Figure S7. GeneEMBED is robust to false negative and false positive edges, Related to Figure 1**. (A) Edges were synthetically and randomly deleted from the Brain network to test sensitivity of GeneEMBED to false negative edges. In blue are plots of precision and recall of GeneEMBED identified genes at various levels of randomly deleted edges. In red are plots of precision and recall of GeneEMBED identified genes when randomly deleted edges are targeted for known (previously identified) genes. (B) Edges were synthetically and randomly added to the Brain network to test sensitivity of GeneEMBED to false positive edges. The plot shows precision and recall of GeneEMBED identified genes at various levels of synthetically added edges. X-axis of '% Edges Added' is relative to the original network size, e.g. at 100%, ~48k edges are randomly added.

7