

Supplementary Materials of “Prompt-based Generation of Natural Language Explanations of Synthetic Lethality for Cancer Drug Discovery”

Ke Zhang^{1,2,3}, Yimiao Feng^{1,4}, Jie Zheng^{1,5,*}

¹ School of Information Science and Technology, ShanghaiTech University, Shanghai, China

²Shanghai Institute of Microsystem and Information Technology, Chinese Academy of Sciences, Shanghai, China

³University of Chinese Academy of Sciences, Beijing, China

⁴Lingang Laboratory, Shanghai, China

⁵Shanghai Engineering Research Center of Intelligent Vision and Imaging, ShanghaiTech University, Shanghai, China

{zhangke1, fengym, zhengjie}@shanghaitech.edu.cn

1 An example of data augmentation prompt

You are a helpful assistant that rephrases text and makes sentences smooth. I will give you a sample, please rephrase the partial sentence after the word “because” of the sample, then give me 10 rephrased answers. Each answer should include the exact noun phrases which I will give you, and each answer must start with “because”. The complete sample is: *TP53 and CDK2 have a synthetic lethality relationship, because TP53 is a tumor suppressor that regulates cell cycle arrest, apoptosis and DNA repair, and CDK2 is a cyclin-dependent kinase that controls cell cycle progression and DNA replication. Therefore, inhibition of CDK2 in TP53-mutant cells results in synergistic cell death due to impaired DNA repair and increased DNA damage.* The phrases are “DNA repair”, “DNA damage”, “cell cycle progression”, “cell death”.

2 Human annotation pipeline

Algorithm 1: Human annotation

```
1 Input: GenepairsCollection  $Q_{nb}$ , AnswerCollection  $M_{nb}$ , CitationCollection  $R_{nb}$ 
2 FactGenePairsCollection  $Q_{fact} \leftarrow \{\}$ 
3 FactAnswerCollection  $M_{fact} \leftarrow \{\}$ 
4 FactCitationCollection  $R_{fact} \leftarrow \{\}$ 
5 FeatureCollection  $F_{fact} \leftarrow \{\}$ 
6 HypotheticalAnswerCollection  $M_{hypo} \leftarrow \{\}$ 
7 for  $i \leftarrow 1$  to  $len(M_{nb})$  do
8   | gene pair  $(u, v) = Q_{nb}[j]$ 
9   | answer  $\leftarrow M_{nb}[i]$ 
10  | citations  $\leftarrow R_{nb}[i]$ 
11  | if AnnotatorReadandCheck(answer, citations) then
12  |   | features  $\leftarrow$  FeatureAnnotation(answer)
13  |   |  $Q_{fact} \leftarrow (u, v)$ 
14  |   |  $M_{fact} \leftarrow$  answer
15  |   |  $R_{fact} \leftarrow$  citations
16  |   |  $F_{fact} \leftarrow$  features
17  | else
18  |   | add answers to  $M_{hypo}$ 
19  | end
20  |  $pairs_{new} \leftarrow$  AnnotatorMiningNewpairs(citations)
21  | if  $pairs_{new} \neq \emptyset$  then
22  |   |  $answers_{new} \leftarrow$  AnnotatorSummarization(citations)
23  |   | features  $\leftarrow$  FeatureAnnotation(answers)
24  |   |  $Q_{fact} \leftarrow pairs_{new}$ 
25  |   |  $M_{fact} \leftarrow answers_{new}$ 
26  |   |  $R_{fact} \leftarrow$  citations
27  |   |  $F_{fact} \leftarrow$  features
28  | end
29 end
```

3 From a KG subgraph to a personalized KG prompt

According to the KG subgraph in Fig. 1, since RAD52 shares two functions with BRAC2's SL partner genes, we assume that RAD52 and BRAC2 also share the functions. Therefore, the KG prompt for BRAC2 and RAD52 is: BRAC2 and RAD52 may share common functions, including DNA Damage Response, DNA repair.

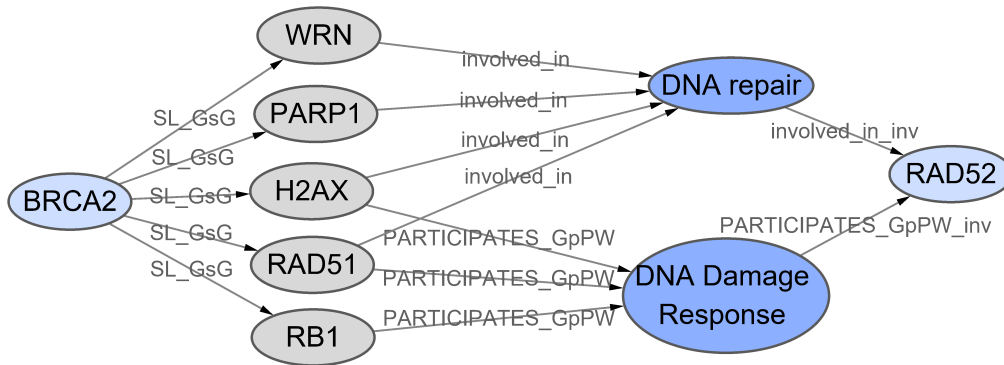


Figure 1: A KG subgraph for an SL gene pair (BRCA2, RAD52). Light blue nodes are the two target genes, grey nodes represent other genes that have SL relationships with BRAC2, and dark blue nodes are two gene functions.