



25 1. INTRODUCTION

26  
27 One of the most critical planning activities in the sciences is  
28 identifying credible priorities for investment. The most high-profile process  
29 of scientific prioritization is the National Academies' Decadal Surveys.

30  
31 A principal challenge faced by this process is the Survey panelists' need to  
32 assess a large -- and rapidly growing -- amount of relevant information,  
33 specifically many tens of thousands of published research papers. The  
34 potential input materials have increased greatly over the years in both  
35 variety and quantity, while the basic processes of the Surveys -- and other  
36 strategic planning activities -- have changed relatively little. The primary  
37 approach for the Surveys over the past half-century (Dressler 2016) remains  
38 the same: a central steering committee of a couple dozen members supported by  
39 large specialty panels. This leads to the primary motivation of our work: are  
40 there ways to substantially improve the current process of identifying the  
41 highest-priority science without adding many additional personnel?

42  
43 We believe that it is time to take advantage of Machine Learning (ML) to  
44 augment the daunting task of determining trends and priorities in science from  
45 a vast amount of information. Advances in ML over the past decade have been  
46 impressive; increasingly powerful ML techniques can comb through a large  
47 corpus of unstructured text to reveal insight into their contents. There have  
48 recently been examples relevant to the process of science prioritization. For  
49 example, Zelnio (2020) reports the successful use of ML to evaluate research  
50 literature for promising technologies, and Krenn et al (2019) demonstrate a  
51 method used to predict future trends in quantum physics.

52  
53 2. METHODOLOGY AND ANALYSIS

54  
55 Our goal was to explore whether an AI-based approach would be able to  
56 significantly enhance human decision-making with regards to high-impact  
57 science research topics. We have created ML models trained on the corpus of  
58 scientific research available in advance of the 2010 Decadal Survey. Our  
59 process is described more fully in Thronson et al. (2021) and Thomas et al.  
60 (2021).

61  
62 Briefly, Natural Language Processing (NLP) was used to process abstracts and  
63 titles drawn from peer-reviewed papers published in the top 10 high-impact  
64 journals in astronomy identified per Thomas (2021) during the time period 1998  
65 to 2010. Papers with abstracts of fewer than 100 characters were filtered out  
66 of the dataset, leaving ~85,000 abstracts. We utilized NLP to extract

67 scientific terms from the abstracts and use these as features in an algorithm  
68 based on Latent Dirichlet Allocation which groups them into research topics.

69  
70 Measurements of the growth in and relative popularity of the topics may be  
71 derived. To determine popularity, we add the fractional contributions that  
72 each topic makes to every abstract in the corpus, or "counts". To determine  
73 growth rate, we calculate the time series of counts for each topic. These time  
74 series may then be analyzed to determine the Compound Annual Growth Rate  
75 (CAGR). The "Research Interest" (RI); that is how much overall interest the  
76 research community places on a given topic, may be quantified from these  
77 measures via

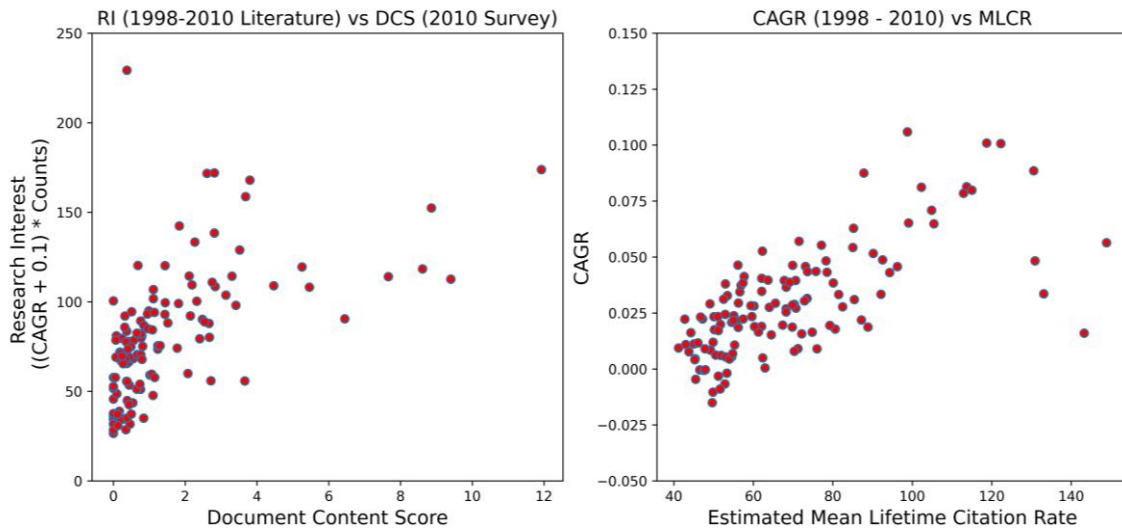
$$78 \quad \quad \quad \text{RI}(t) = (\text{CAGR}(t) + 0.1) * \text{counts}(t) \quad (1)$$

80  
81 where  $t$  is the topic.

82  
83 We next applied these derived topic models to the science frontier panel  
84 chapters 1 - 4 for the 2010 Astronomy and Astrophysics Decadal Survey ("2010  
85 corpus"). After cleaning and extracting features as before, we derived counts  
86 by topic for each of the paragraphs in the 2010 Survey. Paragraphs where the  
87 top three topics contributed less than half of the summed counts total for the  
88 were dropped (i.e., these paragraphs did not have good representation by any  
89 topic models). A "document content score" (or DCS) for each topic was then  
90 derived by taking the remaining paragraphs and summing their counts by topic.

91  
92 These topic models and associated metrics provide a means to quantify and  
93 compare topic content in the literature and the Decadal Survey. We may use  
94 this to check for any common relationships and as a check on the validity of  
95 these metrics. Figures 1a and 1b show the results of this evaluation.

96



97  
 98 **Figure 1a** (left) The 1998 - 2010 literature RI versus the 2010 Survey  
 99 DCS by topic (red dots) indicates a significant, but moderate  
 100 correlation exists. Conversely, only topic CAGR is correlated with the  
 101 estimated topic Mean Lifetime Citation Rate (**Figure 1b**, right).  
 102

103 Figure 1a compares the RI of published abstracts with the 2010 corpus DCS. We  
 104 find a highly significant ( $P < 0.000001$ ) correlation of moderate strength ( $R \cong$   
 105  $0.6$ ), which indicates that research which is both growing in interest and/or  
 106 already has significant research interest is well-represented in the Decadal  
 107 Survey. A separate analysis of the submitted whitepapers (not shown) also  
 108 indicates a similar correlation between RI and the content score of the  
 109 submitted whitepapers.  
 110

111 An essential assumption we have made is that the published body of research  
 112 accurately reflects the interests and priorities of the community of active  
 113 astronomers. In order to help ascertain the validity of this assumption we  
 114 have compared counts, CAGR, and RI for our corpus against the estimated MLCR  
 115 (Thomas 2021) for these same papers as grouped into each topic. Only CAGR was  
 116 found to be correlated with the MLCR ( $R \cong 0.7$ ,  $P < 0.000001$ ; see Figure 1b).  
 117

### 118 3. DISCUSSION

119  
 120 Unlike the MLCR, a metric based on citation rates which are a lagging measure  
 121 of interest in topics of research, these new measures are leading indicators,  
 122 which makes them attractive for use in planning. There appears to be good,  
 123 albeit not perfect, correspondence between the frequency of mention of future  
 124 high-priority research reported in the 2010 Survey and with the content of

125 submitted whitepapers to the RI as determined by the literature of the prior  
126 decade. Interestingly, we find only CAGR to be significantly correlated with  
127 the estimated MLCR. This result suggests that the Decadal Survey places  
128 significant emphasis on established research and may under-emphasize new,  
129 growing research topic areas.

130

131 We note that in all cases our correlations, although significant, are of only  
132 moderate strength and the resultant coefficient of determination ( $R^2$ ), a  
133 measure of how much of the variability in one variable can be "explained by"  
134 variation in the other, is fairly weak ( $R^2 \sim 0.3 - 0.4$ ). Two reasons may  
135 explain why. First, in cases we may be under-sampling the trend in the topic  
136 time series, which would lead to some variation in measured CAGR. An  
137 alternative issue affects measured DCS: our technique models language present  
138 in scientific abstracts, but this may be significantly different from the  
139 language present in the 2010 Survey corpus and could sometimes result in  
140 lowering the DCS values.

141

142 Nevertheless, there is still immediate value in applying this type of analysis  
143 to the science prioritization process. The CAGR measure of topics may be  
144 exploited to identify probable future impactful research topics and papers,  
145 thus creating valuable curated reading. We plan to further understand the  
146 variation and uncertainties in Figure 1, which may make it possible to  
147 distinguish topic regions in these diagrams and provide additional insight.

148

149

150

#### REFERENCES

151

152 Dressler, A., et al. 2016, "The Space Science Decadal Surveys, Lessons Learned  
153 and Best Practices", page 6.

154

155 Krenn, M., and Zelinger, A. 2020, PNAS, 117 (4), 1910-1916  
156 (<https://doi.org/10.1073/pnas.1914370116>)

157

158 Thomas, B., et al. 2021, "Prioritizing Science and Aiding Strategic Planning  
159 in Astronomy by using Artificial Intelligence", AGU Poster Paper, Abstract ID  
160 910791

161

162 Thronson, H., Thomas, B., Barbier, L., & Buonomo, A. 2021, in *Bulletin of the*  
163 *AAS*, 53(1). <https://baas.aas.org/pub/2021n1i541p10>

164

165 Zelnio, R. 2020, private communication.