

## Description of Supplementary Analyses for:

### Mapping Interdisciplinary Fields: Efficiencies, Gaps & Redundancies in HIV/AIDS Research

#### 1 *Topic Models: Determining the Number & Labeling of Topics*

2 The literature has not converged on a single technique in locating topics and topic modelers in fact  
3 appear reluctant to commit to any of the proposed methods. However, several of the developers have  
4 suggested that perplexity scores may be the best option in assessing the generalizability of topic models  
5 (1). Perplexity compares the success of the solutions on the “held out” text. The lower the perplexity  
6 score, the better the model. Blei, Ng, and Jordan (2) describe perplexity as “algebraically equivalent to  
7 the inverse of the geometric mean of per word likelihood” (p.1010). Perplexity is a local maximization  
8 algorithm. As can be seen in Figure S1, for this corpus, it identifies either 30 or 75 topic solutions.

9       Once the model identifies the number of topics, it produces several sets of output. The  
10 probabilistic assignment of papers to each topic are the data used in the primary analyses presented in  
11 the manuscript. In addition to this, the model also provides a list of the top word (stems) associated  
12 with each topic. For our models, we provided the top 50 words associated with each topic to 4 separate  
13 HIV experts, whose collective training/research/teaching experience covers biology, demography,  
14 epidemiology, genetics, clinical medicine, sociology, vaccine development and virology. Each  
15 independently coded the topics. Those codings were summarized by the first author then returned to  
16 each coder for confirmation. In the two cases where differences arose from the independent codings,  
17 focused discussions helped to resolve conflicts. It is important to note that this coding process is only  
18 providing interpretation, not computation/generation of the topic models. Table S1 provides the  
19 applied labels and short descriptions for each of the identified 30 topics.

20       These topics are not equally present within each article in the corpus. The topic model results  
21 provide a proportional allocation of the probability of each abstract deriving from each of the identified

1 topics. Figure S2 presents 2 different ways to use this information to allocate each abstract in the  
2 corpus to the identified list of topics. First, each paper can be considered as being assigned to a single  
3 topic, which we do simply identifying the largest probability for each paper ("Top" topic allocation).  
4 The second option considers that most research likely addressed more than one topic at a time and  
5 thus assigns abstracts to topics proportionally. As can be seen, the distributions do not differ  
6 appreciably depending on which of these allocation strategies we use. In the analyses presented in the  
7 manuscript, all topic allocations (and correspondence comparisons) are conducted using the  
8 proportional allocation strategy. Results did not differ in any appreciable ways if instead using the top  
9 topic allocation approach.

10 The following figures provide several expanded details of the correspondence comparisons  
11 between the identified clustering and subject headings/topics. Figure S3 adds the "consolidated" topics  
12 that are excluded from Figure 2 in the manuscript. Figures S4 & S5 provide the same information  
13 separately for each of the 5-year moving time slices across the observed window, which provide the  
14 basis for several of the interpretations in the discussion section.