1   We thank all the reviewers (**R1**, **R2**, and **R4**) for their thoughtful comments, and respond to questions below.

2   **Hiring as a motivating application & multiple arm pulls per stage (R1, R4):** We were initially motivated by the
3 graduate admissions system run at our university. Here, at every stage, it is possible for *multiple* independent reviewers
4 to look at an applicant. Indeed, our admissions committee strives to hit at least two written reviews per application
5 package, before potentially considering one or more Skype/Hangouts calls with a potential applicant. (In our data, for
6 instance, some applicants received up to 6 independent reviews per stage.)

7   While motivated by academic admissions, we believe our model is of broad interest to industry as well. For example, in
8 tech, it is common to allocate more (or fewer) 30-minute one-on-one interviews on a visit day, and/or multiple pre-visit
9 programming screening teleconference calls. Similarly, in management consulting (cf. the McKinsey citation in our
10 submission), it is common to repeatedly give independent "case study" interviews to borderline candidates.

11   **Information gain (R1):** Information gain $s$ relates to the confidence of accept/reject from an interview vs review. As
12 stages get more expensive, the estimates of utility become more precise, i.e. the estimate comes with a distribution
13 with a lower variance. In practice, a resume review may make a candidate seem much stronger than they are, or a
14 badly written resume could severely underestimate their abilities. However, in-person interviews give better estimates.
15 A strong arm pull with information gain $s$ is equivalent to $s$ weak pulls because it is equivalent to pulling from a
16 distribution with a $\sigma/\sqrt{s}$ sub-Gaussian tail, i.e. it is equivalent to getting a (probably) closer estimate.

17   **Hardness (R1, R2, R4):** Hardness is defined in Eq. 3 as the sum of inverse squared gaps. We will be sure to point to
18 Eq. 3 every time hardness is discussed; the example to be added in response to **R2** will help with intuition as well.

19   **Figure 1 (R4):** Fig. 1 is plotted using fixed values for $K_1$, $\delta$, and $\epsilon$. Hardness is defined in Eq. 3 (see above). We did,
20 however, run experiments to see how the figure shifts when increasing/decreasing each of these values (to quantify the
21 intuition given by Thm. 1). In order to save space we described the effects of the variation on the figure. In the camera
22 ready version we can include a full page of varied graphs in the appendix, the data for which already exists.

23   **Exposition clarification for $T$ and budget (R1):** The definition of $T$ changes based on the section of the paper; we
24 will change this. In Thm. 1 and Alg. 1 the reviewer is correct to assume that $T$ is the cost divided by the information
25 gain. In Thm. 2 and Alg. 2, as presently written, $T$ becomes the budget. We will change $T$ to $\tilde{T}$ in Sec. 4 for clarity.

26   **Budget and $K_i$ (R1):** The reviewer is correct in suggesting that choosing budget and $K_i$ simultaneously is difficult.
27 The policy maker should look at past decisions to estimate gap scores (Eq. 2) and hardness (Eq. 3). There is a clear
28 trade-off between information gain and cost. If the policy maker assumes (based on past data) that the gap scores will be
29 high (it is easy to differentiate between applicants) then the lower stages should have a high $K_i$, and a budget to match
30 the relevant cost $j_i$. If the gap scores are all low (it is hard to differentiate between applicants) then more decisions
31 should be made in the higher, more expensive stages. By looking at the ratio of small gap scores to high gap scores, or
32 by bucketing gap scores, a policy maker will be able to set each $K_i$.

33   **Oracle and suboptimality (R4):** The algorithm assumes that the oracle is correct (which is a very light assumption for
34 the theory section, which assumes a linear objective, but less trivial for general monotone submodular objectives like
35 the diversity-promoting objective used in some of the experiments). However, we find that even with this assumption
36 the algorithm performs well empirically under the submodular function given.

37   **Correctness of Alg. 2 (R1):** Thanks for pointing this out! It is a typo. $\tilde{K}_i$ should be $n - \sum_{a=0}^{i-1} \tilde{K}_a$, i.e., the number of
38 arms remaining. In our code we implemented $\tilde{K}_i$ in that correct way. (See grid test on line 482 in `cut_neurips.py`).

39   **Cost (R1):** The confidence radius was chosen due to a Hoeffding bound. See Appendix B1 for further information.

40   **BRUTaS and CACO differences (R2):** The CACO algorithm is a fixed confidence algorithm, meaning that each stage
41 ends after some (user-supplied) confidence interval is satisfied. When a stage ends, only $K_i$ *arms* move on to the next
42 stage. The BRUTaS algorithm is a fixed budget algorithm where each stage uses a certain amount of *budget*. All arms
43 are pulled $\tilde{T}_{i,t}$ times before a decision to keep or reject a single arm is made. An apples-to-apples example comparing
44 the two algorithms is a bit tricky due to taking different inputs (per-round budget vs. number of arms); we will add two
45 "similar" illustrative examples in the 3-arms/2-stage case to help compare and contrast the two algorithms, though.

46   **Greedy subset selection (R4):** We will run against the greedy algorithm, as suggested. We find the most relevant
47 comparison to be the previous work ("SWAP") from Schumann et. al. 2019. Our method has theoretical guarantees, as
48 does theirs. (Note that we also compare to the actual admissions decision process, which is fairly complex.)

49   **Anonymized data release (R4):** Unfortunately, we are prohibited via FERPA regulations from releasing the real
50 dataset. That said, we are working (i) with our IRB for an anonymized release option as suggested (e.g., by fitting
51 a blanket parameterized model) and (ii) to release through our multi-university research consortium a general set of
52 scripts that mimic our application scraping (e.g., running OCR on SoPs and rec. letters, parsing structured data such as
53 GRE/GPA) and feed that generic data into our algorithms run on their servers. This is a longer process that will require
54 rounds of validation, but our long-term goal is to collect aggregate statistics and push practical impact in this way.