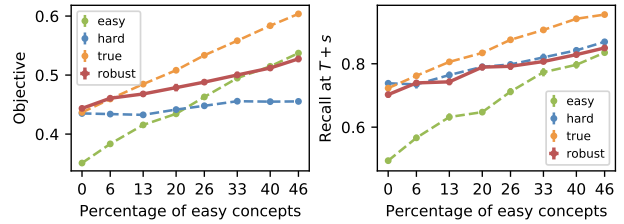


1 We thank the reviewers for their valuable suggestions. Please find our answers (**A**) for each reviewer (**R**) below.

2 **R1, R3: Parameter sensitivity study (In 264–270)**

3 **A:** We had conducted a sensitivity study on simulated
4 learners before choosing the HLR parameters for our
5 user study. These detailed results were omitted from the
6 original submission. We report them below and we will
7 include these results in the revision. In this experiment,
8 we consider two groups of concepts: “easy/common” con-
9 cepts with $\theta = (10, 5, 0)$, and “hard/rare” concepts with
10 $\theta = (3, 1.5, 0)$. Other configurations are kept the same as our user study, with $T = 40$, $n = 15$ and $s = 10$. We vary
11 the number of “easy” concepts from $\{0, 1, \dots, 8\}$ (i.e., up to 50% of the concepts being easy), and consider four types
12 of teachers: (i) “easy”: $\theta = (10, 5, 0)$ for all concepts; (ii) “hard”: $\theta = (3, 1.5, 0)$ for all concepts; (iii) “true”: using
13 the true parameters for each concept; (iv) “robust”: $\theta = (6, 2, 0)$ for all concepts. We plot the performances of these
14 different teachers measured by the two metrics considered in simulations (i.e., the objective value and future recall). As
15 shown in the figures, the “robust” teacher performs well on both metrics, and hence is used for our user study on the
16 German dataset.
17



18 **R1, R3: Time scale of real-life learning settings**

19 **A:** After the publication of this work, it is quite conceivable to apply these ideas in real-life language learning scenarios.
20 We consider collaborating with existing language learning platforms as a natural step for future work.

21 **R1: Fit the parameters of the HLR model, and compare them with the current parameters of choice**

22 **A:** Thanks for the suggestion. Indeed, one can infer θ for each concept from historical data. We have collected 800 user
23 entries from the random teacher on the German dataset (and 3200 entries on Biodiversity), and it is possible to take
24 existing user study histories and fit an HLR model to get an estimate of θ . We plan to include the results in the revision.

25 **R2: Relevant pre-existing work: optimal teaching with exemplars; references on HLR memory model**

26 **A:** Thanks for pointing us to these references. We will certainly include them in the revision. However, there might be
27 some misunderstanding about the differences between the exemplar-based setting (Patil et al. 2011, Nosofsky et al.
28 2018) and the setting of our work. Patil et al. (2011) and Nosofsky et al. (2018) investigated the problem of choosing
29 the optimal exemplars (based on the Generalized Context Model) for teaching a *classification* task; whereas for our case,
30 the exemplars for each class are already given (in other words, we have only one “exemplar” per class), and we aim at
31 optimally teaching the learner to *memorize* the (label of) exemplars. It is unclear how one can adapt the algorithm to
32 our setting, as they are addressing two orthogonal problems. We will explain these points in the updated paper.

33 **R2: Alternative baselines: (1) Optimal “forgetless” learner; (2) different levels of forgetting**

34 **A:** Under our problem setting (as explained in our previous response), if the learner is “forgetless”, then after teaching
35 each concept the recall probability becomes 1, leading to a trivial teaching scenario (by showing each concept once). To
36 see the effect of different levels of forgetting, please refer to our sensitivity study results in response to **R1**.

37 **R2: Significance tests**

38 **A:** We performed χ^2 tests (with contingency tables where rows are algorithms and columns are observed outcomes) and
obtained very similar statistics with Table 1 (see below). We will include these statistics in the updated paper.

	German				Biodiversity				Biodiversity (common)				Biodiversity (rare)			
	GR	LR	RR	RD	GR	LR	RR	RD	GR	LR	RR	RD	GR	LR	RR	RD
p-value (χ^2 tests)	-	0.0652	0.0197	0.0151	-	0.0017	<0.0001	<0.0001	-	0.3111	0.8478	0.0047	-	0.0001	<0.0001	<0.0001

39 **R3: Upper-bound on $F(\pi^g)$**

40 **A:** This is a very interesting suggestion. It will be interesting to establish an upper bound for the greedy or optimal
41 algorithm under particular model configurations, e.g., to provide a necessary condition for achieving a certain target
42 utility under the HLR model (similar to Thm 5). We will further explore this question as future work.
43

44 **R3: Teaching interface: Special consideration for teaching “two consecutive time steps”**

45 **A:** In our teaching interface, there was no gap between two teaching iterations. The learner could copy the answer from
46 the previous iteration to the next question if the same concept was shown. Therefore, we treat this case specially to
47 mitigate such effect. An alternative way is to introduce a small break between two teaching iterations.

48 **R3: Minor Comments**

49 **A:** Thanks for the detailed suggestions. We will fix all the issues with careful proofreading, especially in the Appendix.
50 A few specific answers: (i) *Proof of Thm 1 (Page 14, line 449)*: Yes, we will clarify that $g_i(\cdot) \leq 1$ is due to μ being
51 submodular, and (ii) *Proof of Thm 2 (Page 16)*: Inequality (17) does not require an expectation of the second summand
52 as $(\sigma_{1:t}^g, y_{1:t}^g)$ is the *observed* history. For the inequality between lines 481-482: this is a good point – we will revise this
53 inequality by imposing the condition only on the first part of the expectation. This does not affect the rest of the proof.