REPRESENTATIONAL ALIGNMENT SUPPORTS EFFEC TIVE TEACHING

Anonymous authors

Paper under double-blind review

Abstract

A good teacher should not only be knowledgeable, but should also be able to communicate in a way that the student understands – to share the student's representation of the world. In this work, we introduce a new controlled experimental setting, GRADE, to study pedagogy and representational alignment. We use GRADE through a series of machine-machine and machine-human teaching experiments to characterize a utility curve defining a relationship between representational alignment, teacher expertise, and student learning outcomes. We find that improved representational alignment with a student improves student learning outcomes (i.e., task accuracy), but that this effect is moderated by the size and representational diversity of the class being taught. We use these insights to design a preliminary classroom matching procedure, GRADE-Match, that optimizes the assignment of students to teachers. When designing machine teachers, our results suggest that it is important to focus not only on accuracy, but also on representational alignment with human learners.

023 024 025

026 027

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

The proliferation of digital education resources and AI systems has enabled human and machine 028 teachers to reach potentially millions of students. For example, Massive Open Online Courses 029 (MOOCs) promised to revolutionize education by having top educators record lectures in their domain of expertise and make course materials widely accessible online for many learners. However, 031 this expert-first approach to online learning was not as effective and accessible as hoped for (Reich & Ruipérez-Valiente, 2019), with courses delivered by local teachers often showing better outcomes, 033 in-person and online (Kelly, 2014). More recently, AI systems like ChatGPT have gained hundreds 034 of millions of users, many of whom are using them, or the educational applications they power, to learn new subjects. While these systems can now outperform humans on some tasks (Strachan et al., 035 2023; Van Veen et al.; Thirunavukarasu et al., 2024), their internal representations are not often human-like (Fel et al., 2022; Muttenthaler et al., 2022), highlighting the distinction between domain 037 expertise and the ability to map knowledge into human-understandable spaces. This tension is neither new nor unique to AI; professors can also be experts in their fields that struggle to communicate knowledge to students (Carter et al., 1987; Hinds et al., 2001). Yet many recent public education 040 proposals explicitly focus on increasing teachers' domain expertise (e.g., Ontario, 2024), and much 041 AI research continues to focus on improving the expertise of the agents being developed. Under-042 standing further factors of effective teaching can help determine strategies for improving outcomes 043 in classrooms.

044 We aim to bring together ideas from the burgeoning subfield of *representational alignment* (Sucholutsky et al., 2023b), machine teaching, and the cognitive science of pedagogy to shed light 046 on further improvements for classrooms. We propose that 1) representational alignment between 047 teachers and students, and 2) the size and diversity of the classroom, are both critical for deter-048 mining the effectiveness of human and machine teaching (Figure 1B). To test this hypothesis, we design a simple modular student-teacher cognitive task environment called "Grid Manipulation of Representational Alignment and Domain Expertise" (GRADE) that enables the experimenter to 051 independently control the teacher's expertise on the task, and the degree to which their representations of the task are similar to the student's (Figure 1A). Through simulations and a study where 052 machines teach humans, we establish the relationship between teacher expertise, teacher-student alignment, and student performance. We find that representationally aligned teachers with a high

error rate on the underlying task can outperform highly accurate but representationally misaligned teachers (Figure 1F). These results suggest that if a teacher adapts their representations to match the student, then the student's learning outcomes can significantly improve. We then extend our task from a teacher interacting with individual students to interacting with a class of representationally diverse students where the material they present is broadcasted to all students in the class (e.g., a lecture; see Figure 1D) and determine that class size and representational diversity moderate the effect of representational alignment on student outcomes (Figure 1G). Finally, we design a preliminary classroom matching procedure, GRADE-Match, that takes into account representational alignment, teacher expertise, and class size to optimize learning outcomes when assigning students to teach-ers (Figure 1E). We find that it outperforms both random assignment and MOOC-style assignment (Figure 1H). Our study emphasizes the importance of considering student-teacher representation alignment – not just teacher expertise – in pedagogical settings. This is especially important for designing AI thought partners that can think with us to help us grow (Collins et al., 2024) and tools that help personalize suggestions to individual students (Wang et al., 2024).



Figure 1: Overview. A: GRADE task; teacher and student receive misaligned grids (numbers only represent re-arranged elements, participants do not see them). Teacher is shown all labels (shown as colors) and reveals one per class to the student. Teacher's error rate is controlled by mislabeling their grid. B: Our hypothesized causal model. C: Dyadic interaction between a teacher and a student. "Student-centric" teachers infer the student's representations, making them fully aligned. D: "Classroom" setting where teacher broadcasts examples to all students (who have individual differences in representations); student-centric teachers jointly optimize over all students in the class. E: "School" setting where teachers are matched with students; each student is matched with a single teacher. F: Utility curve relating teacher error rate, representational alignment, and student accuracy in simulations. G: Average accuracy and standard errors in a student-centric class as a function of class size in simulations. H: Average accuracy and standard errors across a school achieved by different matching procedures in simulations.

108
109
1102GRADE: GRID MANIPULATION OF REPRESENTATIONAL ALIGNMENT
AND DOMAIN EXPERTISE

We designed a new controlled task domain, called GRADE, where stimuli are arranged on a grid 111 and labeled. Teachers know all the labels for the stimuli, though their expertise can be modulated by 112 corrupting the true labels to get teachers with varying accuracy. Students do not yet know the labels 113 but do see the stimuli. We focus on the setting where the teacher selects some labeled examples to 114 reveal to the student. The arrangement of the stimuli on the grid might vary between the student and 115 teacher, allowing us to manipulate representational alignment. We show an example of misaligned 116 grids with two labeled examples chosen by the teacher in Figure 1A. GRADE lets us use any stim-117 uli that can be arranged on a grid (e.g., based on pre-set features, as in the "salient-dinos" case in 118 our human experiments, or even amortized embeddings). Here, we define representation alignment with respect to stimuli locations on the student and teacher grids; that is, we compute the Euclidean 119 distance between pairwise swaps between stimuli. GRADE permits modularly-specifiable represen-120 tation functions; we refer to Sucholutsky & Griffiths for a survey of a myriad of ways of measuring 121 representation alignment. Additionally, we focus on a one-shot case where the teacher makes one 122 round of selections. However, researchers can easily extend GRADE to multi-turn interactions. 123 Appendix C contains a theoretical formalism of our setting. 124

125 2.1 RQ1: Does representational alignment affect teaching outcomes?

We begin to explore the relationship between representation alignment and student outcomes in two settings. We instantiate GRADE with two kinds of stimuli that can be arranged on an $N \times N$ grid with K underlying classes: **simple-features** (where each stimulus is only represented by its (x, y) coordinates; see Figure 7) and **salient-dinos** (stick figure images with features varying based on grid location). First, we simulate student-teacher interactions with simple 1-NN agents in the simple-features setting. We then generalize these findings with real human learners in both tasks. We include details of our teacher and student models, and our human experiment, in Appendix E.

Representational (mis)alignment in simulations. In Figure 1F, we trace out a relationship between student-teacher representational alignment, teacher error rate, and student accuracy. We uncover instances wherein students can achieve higher performance by learning from teachers who are *more erroneous* ("less expert") provided the teachers are representationally aligned with the students than comparatively more expert but misaligned teachers, underscoring that it is not just the accuracy of a teacher that matters for student learning outcomes. For a fixed teacher error rate, higher representational alignment is always better for a student (provided the error rate is not too high). We uncover similar curves across grid sizes and the number of categories (see Appendix E).

Representational alignment of machine teachers and human students. We then generalize our 140 findings through human experiments with N = 480 participants (see Appendix E). We construct 141 a utility curve paralleling our simulations by post-hoc varying teacher error rate (see Appendix G). 142 We find in Figure 2 that across both tasks (simple-features and salient-dinos), generally, higher rep-143 resentational alignment induces higher average student accuracy, and report correlations in Table 1. 144 We find that even large increases in teacher error rate can be offset by increasing representational 145 alignment (e.g., a teacher with error rate 0 and representational alignment of 0.3, has similar stu-146 dent outcomes as a teacher with error rate 0.4 and representational alignment 0.8). However, we 147 note that the ordering of high representational alignment is less clear, particularly for the settings 148 where each class corresponds to a column. We posit that people have a strong prior against classes 149 being distributed as columns, and find that especially for the column conditions, participants would 150 often label using strategies that did not correspond to nearest neighbor classification (e.g., several participants labeled in a way that corresponded to different types of tilings). 151

152 153

2.2 RQ2: How does class size and student representational diversity

MODERATE THE EFFECT OF REPRESENTATIONAL ALIGNMENT ON STUDENT OUTCOMES?

154 We have demonstrated that both a teacher's representational alignment and their accuracy matter 155 for student outcomes. So far, our teachers have been self-centered; they use a single representa-156 tion to select labeled examples. This approach suits machine teachers unable to adapt to specific 157 students but does not capture capabilities of adaptable human or future machine teachers. Addi-158 tionally, classroom teachers often address multiple students with differing representations, making 159 example selection more complex. Here, we consider student-centric teachers who aim to maximize the average performance of a student pool by simulating likely student learning outcomes to 160 various selections in an "inner loop" optimization (see Appendix F). While our earlier findings 161 suggest student-centric teachers may enhance learning by becoming representationally aligned with



Figure 2: Relating teacher error rate and representational alignment between machine teachers and human students to student accuracy. From left to right: *simple-features* one class per quadrant; *salient-dinos* one per quadrant; *simple-features* one class per column (6); *salient-dinos* one class per column (7). Format follows Figure 1F.

students, we hypothesize that this effect will be moderated by the group's size and representational diversity, since teachers must optimize for all students simultaneously.

Setup We extend GRADE to investigate classrooms of varying sizes. Because we sample students for each class from the same pool of representationally diverse students (Appendix F), increasing classroom size will generally increase representational diversity.

Results We investigate the relationship between classroom size and student performance by sampling teachers with a range of error rates (between 0 and 0.5 in increments of 0.1) and classroom sizes (10 seeds per setting). Each student-centered teacher optimizes for their class through T = 100inner loop iterations. We then marginalize over our sampled error rates to compute an expected average classroom accuracy per classroom size. We find that the performance of students in a classroom with a student-centric teacher is initially high (i.e., with a class of only a single student, the studentcentric teacher would be equivalent to a fully representationally aligned teacher in the dyad setting) but falls off rapidly as a function of classroom size and then plateaus (as shown in Figure 1G).

185 2.3 RQ3: CAN WE MATCH TEACHERS AND STUDENTS TO IMPROVE OUTCOMES?

Given a "school" of teachers and students, how can we simultaneously group students into "classrooms" and determine which teacher to allot to which class? We begin to explore this question through a series of "classroom matching" experiments. We develop a *classroom matching procedure*, GRADE-Match, which given a pool of students and teachers, assigns groups of students to teachers based on our representational alignment-teacher utility curve. We emphasize though, that our analogy to "classrooms" and "schools" is explored in simulation with machines teaching machines; substantial future work is required to investigate the generalization of possible links between representational alignment, teacher error rate, and classroom properties in practice.

¹⁹³ Setup

168

169

170

Student and teacher pools. We focus on our simple-features setting and extend our dyad (single teacher, single student) setting to simulated pools of teachers and students over our same 6×6 grid. We design two different pools of students and teachers (unstructured and structured). We include pool construction and generalization experiments to the salient-dinos setting in Appendix F.

198 Matching procedures. We propose matching students using our utility curve to estimate their ac-199 curacy (Grade-Match (Ours)). We compute the representational alignment between a student and teacher and index into a bucketed version of the utility curve (recomputed by averaging over sam-200 ples of corrupted students; see Appendix F). that we construct in Section 2 using both the repre-201 sentational alignment and teacher's expected error rate (which we assume we have access to). The 202 resulting metric is the student's expected performance under a specific teacher and classroom. We 203 iterate over all teachers for each student and select the teacher who helps the student achieve the 204 highest expected performance. We consider three baselines: (i) Random matching of students to 205 teachers, (ii) **MOOC** which matches all students to the lowest error rate teacher, and (iii) **Optimal** 206 wherein we use a brute force-search to match students to the highest attainable accuracy, giving an 207 indication of the upper limit of performance that a matching algorithm could possibly achieve. Gaps 208 between (ii) and (iii) further drive home the importance of going beyond teacher accuracy when 209 pairing students to teachers.

Results Our matching algorithm, which groups students to teachers based on their representational alignment and teacher error rate, generally outperforms random matching and, particularly for top-performing students, is better than having assigned the student to an expert (minimal error rate teacher; MOOC) who may be representationally distinct (see Figure 1H and Appendix Tables 2 and 3). This observation is intriguing – students may not achieve their full potential when paired with a representationally misaligned teacher, even if that teacher is an expert. We observe performance gains for our utility curve-based matching across both pool types. However, we do not yet attain optimal matching performance, perhaps due to a mismatch in our utility curve.

216 3 DISCUSSION AND LIMITATIONS

Expertise on a task is not sufficient to be a good teacher; representational alignment matters too. 218 Using a new controlled experimental paradigm (GRADE), we trace out a utility curve between 219 teacher accuracy, teacher-student representational alignment, and student accuracy to characterize 220 the crucial relationship between representational (mis)alignment and student learning outcomes. 221 We put this utility curve to work to better match teachers to students based on representational 222 alignment. Our work underscores the importance of teachers representing a diversity of students 223 and arranging student-teacher groups to ensure there is at least one teacher that any student can 224 effectively learn from. This motivates further investigation into representational alignment and its influence on pedagogical effectiveness in multiple learning settings, like teacher-student interactions 225 and peer mentorship. Yet, we emphasize that our work is a first step in the study of the relationship 226 between representational alignment, teacher efficacy, and student-teacher matching. Our simulations 227 always assume that students are 1-NN classifiers, which grossly undercuts the richness of human 228 behavior. Further, our simulated students' representations are fixed; in practice, students adapt their 229 representations over time. We also only consider single-turn, single-lesson settings, wherein students 230 have no indication of the reliability of the teacher. We look forward to investigations that leverage 231 and extend GRADE to go beyond our simple yet revealing initial setting. 232

233 234 REFERENCES

238

239

240

243

244

245

246

260

261

- Rosie Aboody, Joey Velez-Ginorio, Laurie R Santos, and Julian Jara-Ettinger. When naive pedagogy
 breaks down: Adults rationally decide how to teach, but misrepresent learners' beliefs. *Cognitive Science*, 47(3):e13257, 2023.
 - John R Anderson, C Franklin Boyle, and Brian J Reiser. Intelligent tutoring systems. *Science*, 228 (4698):456–462, 1985.
- Jacob Andreas and Dan Klein. Reasoning about pragmatics with neural listeners and speakers. *arXiv preprint arXiv:1604.00562*, 2016.
 - Ilona Bass, Elizabeth Bonawitz, Daniel Hawthorne-Madell, Wai Keen Vong, Noah D Goodman, and Hyowon Gweon. The effects of information utility and teachers' knowledge on evaluations of under-informative pedagogy across development. *Cognition*, 222:104999, 2022.
- Benjamin S Bloom. The 2 sigma problem: The search for methods of group instruction as effective as one-to-one tutoring. *Educational Researcher*, 13(6):4–16, 1984.
- Elizabeth Bonawitz, Patrick Shafto, Hyowon Gweon, Noah D Goodman, Elizabeth Spelke, and Laura Schulz. The double-edged sword of pedagogy: Instruction limits spontaneous exploration and discovery. *Cognition*, 120(3):322–330, 2011.
- Sophie Bridgers, Julian Jara-Ettinger, and Hyowon Gweon. Young children consider the expected utility of others' learning to decide what to teach. *Nature Human Behaviour*, 4(2):144–152, 2020.
- 255 Kathy Carter, Donna Sabers, Katherine Cushing, Stefinee Pinnegar, and David C. 256 Berliner. Processing and using information about students: A study of expert, 257 novice, and postulant teachers. Teaching and Teacher Education, 3(2):147–157, 258 1987. ISSN 0742-051X. doi: https://doi.org/10.1016/0742-051X(87)90015-1. URL. https://www.sciencedirect.com/science/article/pii/0742051X87900151. 259
 - Alicia M Chen, Andrew Palacci, Natalia Vélez, Robert Hawkins, and Samuel J Gershman. A hierarchical Bayesian approach to adaptive teaching, Dec 2022.
- Katherine M Collins, Ilia Sucholutsky, Umang Bhatt, Kartik Chandra, Lionel Wong, Mina Lee,
 Cedegao E Zhang, Tan Zhi-Xuan, Mark Ho, Vikash Mansinghka, et al. Building machines that
 learn and think with people. *Nature Human Behaviour*, 8(10):1851–1863, 2024.
- Gergely Csibra and György Gergely. Natural pedagogy. *Trends in Cognitive Sciences*, 13(4):148–153, 2009.
- 269 Robert Dale and Ehud Reiter. Computational interpretations of the Gricean maxims in the generation of referring expressions. *Cognitive Science*, 19(2):233–263, 1995.

270	Thomas S Dee. Teachers, race, and student achievement in a randomized experiment.	Review of
271	<i>Economics and Statistics</i> , 86(1):195–210, 2004.	5
272		

- Anna J. Egalite, Brian Kisida, and Marcus A. Winters. Representation in the classroom: The 273 effect of own-race teachers on student achievement. Economics of Education Review, 45: 274 44-52, 2015. ISSN 0272-7757. doi: https://doi.org/10.1016/j.econedurev.2015.01.007. URL 275 https://www.sciencedirect.com/science/article/pii/S0272775715000084. 276
- 277 Thomas Fel, Ivan F Rodriguez Rodriguez, Drew Linsley, and Thomas Serre. Harmonizing the object 278 recognition strategies of deep neural networks with humans. Advances in Neural Information 279 Processing Systems, 35:9432–9446, 2022.
- Christian Fischer, Zachary A. Pardos, Ryan Shaun Baker, Joseph Jay Williams, Padhraic Smyth, 281 Renzhe Yu, Stefan Slater, Rachel Baker, and Mark Warschauer. Mining big data in educa-282 tion: Affordances and challenges. *Review of Research in Education*, 44(1):130–160, 2020. doi: 283 10.3102/0091732X20903304. URL https://doi.org/10.3102/0091732X20903304. 284
- Michael Frank. Modeling the dynamics of classroom education using teaching games. In Proceedings of the Annual Meeting of the Cognitive Science Society, volume 36, 2014. 286
- 287 Daniel Fried, Ronghang Hu, Volkan Cirik, Anna Rohrbach, Jacob Andreas, Louis-Philippe 288 Morency, Taylor Berg-Kirkpatrick, Kate Saenko, Dan Klein, and Trevor Darrell. Speaker-follower 289 models for vision-and-language navigation. Advances in Neural Information Processing Systems, 31, 2018. 291
- Noah D Goodman and Michael C Frank. Pragmatic language interpretation as probabilistic infer-292 ence. Trends in Cognitive Sciences, 20(11):818-829, 2016. 293
- 294 Herbert P Grice. Logic and conversation. In Speech Acts, pp. 41–58. Brill, 1975. 295

- 296 Hyowon Gweon. Inferential social learning: Cognitive foundations of human social learning and 297 teaching. Trends in Cognitive Sciences, 25(10):896–910, 2021.
- 298 G Harfitt. The role of the community in teacher preparation: Exploring a different pathway to 299 becoming a teacher. Front. Educ. 3: 64. doi: 10.3389/feduc, 2018. 300
- 301 Pamela J Hinds, Michael Patterson, and Jeffrey Pfeffer. Bothered by abstraction: the effect of exper-302 tise on knowledge transfer and subsequent novice performance. Journal of Applied Psychology, 303 86(6):1232, 2001.
- Mark K Ho, Michael Littman, James MacGlashan, Fiery Cushman, and Joseph L Austerweil. Show-305 ing versus doing: Teaching by demonstration. Advances in Neural Information Processing Sys-306 tems, 29, 2016.
- 308 Dun-Ming Huang, Pol Van Rijn, Ilia Sucholutsky, Raja Marjieh, and Nori Jacoby. Characterizing similarities and divergences in conversational tones in humans and llms by sampling with people. arXiv preprint arXiv:2406.04278, 2024. 310
- 311 Andrew P Kelly. Disruptor, distracter, or what? a policymaker's guide to massive open online 312 courses (MOOCs). Bellwether Education Partners, 2014. 313
- 314 Michalis Korakakis and Andreas Vlachos. Improving the robustness of NLI models with minimax training. In Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki (eds.), Proceedings 315 of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long 316 Papers), ACL 2023, pp. 14322-14339, 2023. doi: 10.18653/V1/2023.ACL-LONG.801. URL 317 https://doi.org/10.18653/v1/2023.acl-long.801. 318
- 319 Laurent Lessard, Xuezhou Zhang, and Xiaojin Zhu. An optimal control approach to sequential 320 machine teaching. In The 22nd International Conference on Artificial Intelligence and Statistics, 321 pp. 2495–2503, 2019. 322
- Jessy Lin, Daniel Fried, Dan Klein, and Anca Dragan. Inferring rewards from language in context. 323 arXiv preprint arXiv:2204.02515, 2022.

333

351

352

353

- Ulf Liszkowski, Malinda Carpenter, and Michael Tomasello. Twelve-month-olds communicate helpfully and appropriately for knowledgeable and ignorant partners. *Cognition*, 108(3):732– 739, 2008.
- Ji Liu and Xiaojin Zhu. The teaching dimension of linear learners. *Journal of Machine Learning Research*, 17(162):1–25, 2016.
- Ryan Liu, Jiayi Geng, Joshua C Peterson, Ilia Sucholutsky, and Thomas L Griffiths. Large language
 models assume people are more rational than we really are. *arXiv preprint arXiv:2406.17055*, 2024.
- Weiyang Liu, Bo Dai, Ahmad Humayun, Charlene Tay, Chen Yu, Linda B Smith, James M Rehg, and Le Song. Iterative machine teaching. In *International Conference on Machine Learning*, pp. 2149–2158. PMLR, 2017.
- Weiyang Liu, Bo Dai, Xingguo Li, Zhen Liu, James Rehg, and Le Song. Towards black-box iterative
 machine teaching. In *International Conference on Machine Learning*, pp. 3141–3149, 2018.
- Weiyang Liu, Zhen Liu, Hanchen Wang, Liam Paull, Bernhard Schölkopf, and Adrian Weller. It erative teaching by label synthesis. *Advances in Neural Information Processing Systems*, 34: 21681–21695, 2021.
- Yuzhe Ma, Robert Nowak, Philippe Rigollet, Xuezhou Zhang, and Xiaojin Zhu. Teacher Improves
 Learning by Selecting a Training Subset. In Amos Storkey and Fernando Perez-Cruz (eds.),
 Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics,
 volume 84 of *Proceedings of Machine Learning Research*, pp. 1366–1375, 09–11 Apr 2018. URL
 https://proceedings.mlr.press/v84/ma18a.html.
- Maya Malaviya, Ilia Sucholutsky, Kerem Oktar, and Thomas L Griffiths. Can humans do less-than one-shot learning? In 44th Annual Meeting of the Cognitive Science Society: Cognitive Diversity,
 CogSci 2022, 2022.
 - Smitha Milli and Anca D Dragan. Literal or pedagogic human? analyzing human model misspecification in objective learning. In *Uncertainty in Artificial Intelligence*, pp. 925–934. PMLR, 2020.
- Lukas Muttenthaler, Jonas Dippel, Lorenz Linhardt, Robert A Vandermeulen, and Simon Kornblith.
 Human alignment of neural network representations. *arXiv preprint arXiv:2211.01201*, 2022.
- Jakob Niedermann, Ilia Sucholutsky, Raja Marjieh, Elif Celen, Thomas L Griffiths, Nori Jacoby, and Pol van Rijn. Studying the Effect of Globalization on Color Perception using Multilingual Online Recruitment and Large Language Models, Feb 2024. URL osf.io/preprints/psyarxiv/3jvxw.
- Kerem Oktar, Ilia Sucholutsky, Tania Lombrozo, and Thomas L. Griffiths. Dimensions of disagree ment: Unpacking divergence and misalignment in cognitive science and artificial intelligence,
 2023.
- 364 Ontario, May 2024. URL https://news.ontario.ca/en/backgrounder/1004649/modern-relevant-and-365
- Stefan Palan and Christian Schitter. Prolific.ac—A subject pool for online experiments. *Journal of Behavioral and Experimental Finance*, 17:22–27, 2018.
- John F Pane, Beth Ann Griffin, Daniel F McCaffrey, and Rita Karam. Effectiveness of cognitive tutor algebra I at scale. *Educational Evaluation and Policy Analysis*, 36(2):127–144, 2014.
- Zeju Qiu, Weiyang Liu, Tim Z Xiao, Zhen Liu, Umang Bhatt, Yucen Luo, Adrian Weller, and
 Bernhard Schölkopf. Iterative teaching by data hallucination. In *International Conference on Artificial Intelligence and Statistics*, pp. 9892–9913, 2023.
- Sunayana Rane, Mark Ho, Ilia Sucholutsky, and Thomas L Griffiths. Concept alignment as a pre requisite for value alignment. *arXiv preprint arXiv:2310.20059*, 2023.
- 377 Justin Reich. *Failure to disrupt: Why technology alone can't transform education*. Harvard University Press, 2020.

378	Justin Reich and José A Ruipérez-Valiente. The MOOC pivot. Science, 363(6423):130–131, 2019.
379	Jonnifer King Dice. Teacher quality: Understanding the effectiveness of teacher attributes. EDIC
300	2003
382	2003.
302	Patrick Shafto, Noah D Goodman, and Thomas L Griffiths. A rational account of pedagogical
384	reasoning: Teaching by, and learning from, examples. <i>Cognitive Psychology</i> , 71:55–89, 2014.
385	James Strachan, Dalila Albergo, Giulia Borghini, Oriana Pansardi, Eugenio Scaliti, Alessandro
386	Rufo, Guido Manzi, Michael Graziano, and Cristina Becchio. Testing theory of mind in GPT
387	models and humans. 2023.
388	Ilia Sucholutsky and Thomas L Griffiths Alignment with human representations supports robust
389	few-shot learning. <i>NeurIPS</i> , 2023.
390	
391	Ilia Sucholutsky, Ruairidh M Battleday, Katherine M Collins, Raja Marjieh, Joshua Peterson, Pulkit
392	of supervision signals. In Uncertainty in Artificial Intelligence, pp. 2036–2046. PMLR, 2023a
393	of supervision signals. In oncertainty with typetar interaction, pp. 2050–2010. I filler, 2020a.
394	Ilia Sucholutsky, Lukas Muttenthaler, Adrian Weller, Andi Peng, Andreea Bobu, Been Kim,
396	Bradley U. Love, Erin Grant, Jascha Achterberg, Joshua B. Tenenbaum, Katherine M. Collins, Katherine I. Hermann, Kerem Oktor, Klaus Graff, Martin N. Habert, Nori Jacoby, Oliver, Zhang,
397	Raia Marijeh, Robert Geirhos, Sherol Chen, Simon Kornblith, Sunavana Rane, Talia Konkle
398	Thomas P. O'Connell, Thomas Unterthiner, Andrew K. Lampinen, Klaus-Robert Müller, Mariya
399	Toneva, and Thomas L. Griffiths. Getting aligned on representational alignment, 2023b.
400	Theodore Sumare Depart Hawking Mark K He Tom Criffiths and Dylan Hadfald Manall Hey to
401	talk so AI will learn: Instructions descriptions and autonomy Advances in Neural Information
402	Processing Systems, 35:34762–34775, 2022.
403	
404	Theodore R Sumers, Mark K Ho, Robert D Hawkins, Karthik Narasimhan, and Thomas L Griffiths.
405	Intelligence volume 35 pp 6002–6010 2021
406	<i>Intelligence</i> , volume 55, pp. 0002–0010, 2021.
407	Theodore R Sumers, Mark K Ho, Robert D Hawkins, and Thomas L Griffiths. Show or Tell?
409	Exploring when (and why) teaching with language outperforms demonstration. Cognition, 232: 105326, 2023
410	103320, 2023.
411	Arun James Thirunavukarasu, Shathar Mahmood, Andrew Malem, William Paul Foster, Rohan
412	Sanghera, Refaat Hassan, Sean Zhou, Shiao Wei Wong, Yee Ling Wong, Yu Jeat Chong, Ab-
413	Rauz Daniel Shu Wei Ting and Darren Shu Jeng Ting Large language models approach
414	expert-level clinical knowledge and reasoning in ophthalmology: A head-to-head cross-sectional
415	study. PLOS Digital Health, 3(4):1–16, 04 2024. doi: 10.1371/journal.pdig.0000341. URL
416	https://doi.org/10.1371/journal.pdig.0000341.
417	Dave Van Veen, Cara Van Uden, Louis Blankemeier, Jean-Benoit Delbrouck, Asad Aali, Christian
410	Bluethgen, Anuj Pareek, Malgorzata Polacin, Eduardo Pontes Reis, Anna Seehofnerova, et al.
420	Clinical text summarization: Adapting large language models can outperform human experts.
421	Research Square.
422	Natalia Vélez, Alicia M Chen, Taylor Burke, Fiery A Cushman, and Samuel J Gershman. Teachers
423	recruit mentalizing regions to represent learners' beliefs. <i>Proceedings of the National Academy</i>
424	of Sciences, 120(22):e2215015120, 2023.
425	Huanhuan Wang Ahmed Tlili Ronghuai Huang Zhenyu Cai Min Li Zui Cheng Dong Yang
426	Mengti Li, Xixian Zhu, and Cheng Fei. Examining the applications of intelligent tutoring systems
427	in real educational contexts: A systematic literature review from the social experiment perspec-
428	tive. Education and Information Technologies, 28(7):9113–9148, 2023.
430	Pei Wang, Jungi Wang, Pushpi Paranamana, and Patrick Shafto. A mathematical theory of coop-
431	erative communication. Advances in Neural Information Processing Systems, 33:17582–17593,
	2020.

432 433	Rose E Wang, Ana T Ribeiro, Carly D Robinson, Susanna Loeb, and Dora Demszky. Tutor copilot: A human-ai approach for scaling real-time expertise. <i>arXiv preprint arXiv:2410.03017</i> , 2024.
434	Sida I Wang, Percy Liang, and Christopher D Manning. Learning language games through interac-
436	tion. arXiv preprint arXiv:1606.02447, 2016.
437 438	Andrea Wynn, Ilia Sucholutsky, and Thomas L Griffiths. Learning human-like representations to
439	enable learning numan values. arxiv preprint arxiv:2312.14106, 2023.
440	Teresa Yeo, Parameswaran Kamalaruban, Adish Singla, Arpit Merchant, Thibault Asselborn, Louis
441	Faucon, Pierre Dillenbourg, and Volkan Cevher. Iterative classroom teaching. In <i>Proceedings of</i>
442	the AAAI Conjerence on Artificial Intelligence, volume 55, pp. 5084–5092, 2019.
443 444	Chen Zhang, Xiaofeng Cao, Weiyang Liu, Ivor Tsang, and James Kwok. Nonparametric iterative machine teaching. In <i>International Conference on Machine Learning</i> , pp. 40851–40870, 2023.
445	indennie edennig. In International Conference on Indennie Learning, pp. 10051 10010, 2020.
446 447	Tan Zhi-Xuan, Lance Ying, Vikash Mansinghka, and Joshua B. Tenenbaum. Pragmatic instruction following and goal assistance via cooperative language-guided inverse planning, 2024.
448	
449	optimal education. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 29,
450	2015.
451	Xiaojin Zhu, Adish Singla, Sandra Zilles, and Anna N Rafferty. An overview of machine teaching
452	arXiv preprint arXiv:1801.05927, 2018.
453	
454	
455	
456	
457	
458	
459	
400	
462	
463	
464	
465	
466	
467	
468	
469	
470	
471	
472	
473	
474	
475	
476	
477	
478	
479	
480	
481	
482	
483	
484	
485	

486 A BROADER IMPACT AND SOCIETAL RISKS

488 As we discuss in Section 3, our work portends broader implications for the design of machine-human 489 teaching setups where machines are intentionally built with representation alignment in mind, as 490 well as representation diversity to safeguard against threats to inclusivity. It is possible that our 491 simulations could support interventions in real classrooms, e.g., informing classroom size decisions 492 drawing on measures of the representational diversity of a classroom pool and the expertise of the teacher. However, we heed caution in over-generalizing our results to settings where real student 493 494 experiences and learning potential is at stake. Broad-brush application of AI systems in education has not been met with universal success (Reich, 2020) and inappropriately incorporated can have 495 unintended impacts on student success (Fischer et al., 2020). 496

497 498

499

B RELATED WORK

500 Learning Sciences. The extended learning sciences community has studied aspects of what makes 501 for a good teacher or computer-based teaching system. The expertise or quality of the teacher with 502 respect to excellence of schooling, certification, and a teacher's own test scores have been observed 503 to positively affect student learning (Rice, 2003). More classroom-adaptive qualities, like a teacher's amount of experience in classrooms and teaching strategies employed (i.e., pedagogy) are also top 504 contributing attributes (Rice, 2003). Closeness of representation to students with respect to demo-505 graphic features has been shown to lead to more effective student performance (Dee, 2004), in part 506 due to the role model effect, but also because teachers closer in these dimensions can serve as so-507 ciocultural interlocutors, helping translate the relevance of material to students (Egalite et al., 2015; 508 Harfitt, 2018). Intelligent Tutoring Systems (Anderson et al., 1985), growing out of the cognitive and 509 learning sciences, have been a consistently effective paradigm of computer-based teaching (Wang 510 et al., 2023), primarily utilizing the pedagogy of mastery learning (Bloom, 1984). They adapt the 511 amount of prescribed practice based on a representation of the student's level of mastery of the skill 512 being worked on and provide procedural remediation in the problem-solving context. In a two-year, 513 large-scale evaluation, a commercial ITS was found to be effective overall, but only demonstrated 514 superior learning gains to standard classroom instruction in the second year. It was hypothesized 515 that this may have been due to teachers needing to learn how best to align their classroom to the technology (Pane et al., 2014). 516

517 Machine teaching. Machine teaching aims to study the problem of teaching efficiency by charac-518 terizing such efficiency as the minimal number of effective data examples that is needed for a learner 519 to learn some target concept. It has an intrinsic connection to optimal education (Zhu, 2015), cur-520 riculum learning (Liu et al., 2017; Korakakis & Vlachos, 2023) and optimal control (Lessard et al., 2019). Depending on the type of learner, machine teaching can be performed in a batch setting (Zhu, 521 2015; Zhu et al., 2018; Liu & Zhu, 2016) or an iterative setting (Liu et al., 2017; 2018; 2021; Qiu 522 et al., 2023; Zhang et al., 2023). The batch teaching aims to find a training dataset of minimal size 523 such that a learner can learn a target concept based on this minimal dataset. The iterative teach-524 ing seeks a sequence of data such that the learner can sequentially learn the target concept within 525 minimal iterations. Complementary to these works, our findings indicate that, alongside the quality 526 of examples that the teacher selects, it is also critical for both the teacher and the student to share 527 similar representations. 528

Pragmatic communication. Successful communication rests on our ability to understand others' 529 beliefs and intentions (Gweon, 2021; Vélez et al., 2023). Indeed, even young children are sensitive 530 to others' knowledge and competence when teaching (Liszkowski et al., 2008; Bridgers et al., 2020) 531 and learning (Bass et al., 2022; Csibra & Gergely, 2009; Bonawitz et al., 2011) from others. Inspired 532 by Gricean pragmatics (Grice, 1975), recent computational models have formalized this process 533 as recursive reasoning about others' latent mental states (Chen et al., 2022; Goodman & Frank, 534 2016; Shafto et al., 2014). Such pragmatic models have been used to study and facilitate human-AI interaction (Sumers et al., 2021; 2022; Lin et al., 2022; Andreas & Klein, 2016; Dale & Reiter, 1995; 536 Fried et al., 2018; Wang et al., 2016; 2020; Ho et al., 2016; Zhi-Xuan et al., 2024; Liu et al., 2024). 537 Crucially, however, when either party *fails* to accurately model the other's beliefs or perspective, human-human (Aboody et al., 2023; Sumers et al., 2023) and human-AI (Milli & Dragan, 2020; 538 Sumers et al., 2022) communication can be significantly degraded. Our work adds to this literature by formalizing and analyzing the effect of *representational* misalignment on communication.



Figure 3: Schematic of teaching and representational alignment. Teachers and students have distinct representational spaces (X, Y_s) with some mapping between them (T_s) . There is a true label function (f) that can be projected onto both the teacher and student spaces, but a teacher may not perfectly know this true label function and have their own, diverging label function (f'). The teacher designs curricular materials $(L_0;$ a set of examples paired with labels) that are projected to each student's space (L_s) , where each student uses them to learn a label function (g_s) . Each student's performance (V_s) is then measured as the divergence between the learned label function and the hidden true label function $(T'_s(f))$.

564 **Representational alignment.** Representational alignment (Sucholutsky et al., 2023b) offers a con-565 ceptual and grounded mathematical framework for characterizing teaching settings wherein two or 566 more agents engage on some task. Already, ideas from representational alignment are providing 567 new ways of thinking about machine learning efficiency (Sucholutsky et al., 2023a; Sucholutsky 568 & Griffiths, 2023), value alignment (Rane et al., 2023; Wynn et al., 2023), disagreement (Oktar et al., 2023), and applications like human & machine translation and conversation (Niedermann 569 et al., 2024; Huang et al., 2024). In this study, we show that representational alignment is a key 570 dimension in predicting and optimizing student outcomes, with similar importance as the teacher's 571 subject expertise. 572

573 574

575

C THEORETICAL FORMULATION UNDERLYING GRADE

We offer a deeper theoretical formalism for our setting. Figure 3 shows a schematic of our teaching and representation alignment framework. Consider a space X of stimuli. We consider the case where the teacher tries to teach the students some function $f: X \to C$. We illustrate a simple case in Figure 3 wherein C is a binary classification $C = \{0, 1\}$ dividing X into two regions (C = 0 and C = 1 are represented in Fig. 3 in light and dark gray, respectively). The teacher observes label function $f': X \to C$, which may be different from f.

The teacher chooses *n* points from the space $x_1, x_2, \ldots, x_n \in X$ and assigns labels to the points l_i . The teacher materials can be represented by the labeled points : $L_0 = (x_i, l_i)_{i=1,...,n}$. To represent the fact that students' representations may differ from that of the teacher, we assume that the student *s* has a space Y_s that corresponds to the student's representations. Note that each student is part of some classroom or population of students $s \in S$.

Next, we assume there is some transformation $T_s : X \to Y_s$. We assume that the function T_s is also selected from some parametric function $T_s \sim \mathcal{T}$. The student s observes stimuli presented by the teacher $y_i = T_s(x_i)$ and labels l_i . The student learning input (i.e., the teaching materials mapped into the student's space) is thus $L_s = (T_s(x_i), l_i)_{i=1,...,n}$. From that, the student infers the labeling for the rest of the space, which can be represented as the learning function $g_s(y|L_s)$. The classification performance of the student is tested over additional test points where the expected performance of the student is $V_s = D(g_s(y|L_s)||T_s(f(x)))$. Here, D represents some distance measure.

594 D COMPUTE DETAILS

596

597

598

600 601

602

603 604

605

606

607

608

609

610

611

612

613

628

629 630

647

All experiments were run on an 8-core, 16 GB memory laptop. Experiments were run exclusively on CPUs and were all runnable within at most three hours. Our experiments are reproducible and all the implementations for all computational experiments will be made available open-source upon publication.

E ADDITIONAL DETAILS ON TASK SETUP FOR THE SINGLE TEACHER-SINGLE STUDENT SETTING

Student and teacher models We instantiate our student (g_s) as a 1-nearest neighbor (1-NN) classifier, who takes as input the teacher's revealed examples (L_s) and classifies each of the unlabeled points. Student performance (V_s) is computed as the accuracy of their classifications over the unlabeled points. The teacher chooses K' points intended to maximally help the student (whom the teacher "knows" is a 1-NN classifier) to achieve high accuracy on the remaining points. We assume the teacher has access to labels for all cells; however, the "erroneous" teacher with some probability assumes the wrong label on a cell (i.e., f' is different from f, which can ripple into their selections accordingly). The teacher computes the centroid of each class (using its own believed labels f', which may have errors) and selects one example per class to reveal to the student. The teacher reveals its believed labels to the student for the selected points.



Figure 4: Utility curves on a 6×6 grid for different label structures. (Left:) 4 class underlying label-per-quadrant; (**Right:**) 6 class underlying label-per-column.

Constructing the dyadic utility curve To construct our utility curve in Section 2, we sweep over 631 a range of possible teacher error rate parameterizations (from 0 to 0.9 in increments of 0.1) and 632 representation corruption levels (from 0 to 1.0 in increments of 0.01). We always use the same 633 "student" and corrupt teachers over the respective student grid. We compute the representation 634 alignment between the student and corrupted teacher; as the pairwise swaps ("corruptions") are 635 randomly made over a fraction of the grid parameterized by the corruption level, we bucketize 636 the resulting observed representation alignment between student and teacher. We then sample 10 637 different seeds of selections for each teacher and average student performance. We repeat the same 638 sweeps over teacher parameterizations for our two labeling schemes: grids wherein each column 639 corresponds to one label (N labels for an $N \times N$ grid), and one where each quadrant corresponds 640 to one label (four labels). We average the resulting utility curves across label types. 641

Impact of underlying label structure We depict the separate utility curves in Figure 4. Notably, we observe different utility curves for different label structures. While there are some minor rank swaps between teachers across the structures, we see high Spearman rank correlation $(\rho = 0.994, p << 10e - 48)$ between the two settings underscoring general consistency in teacher orderings.

Additional details on experiments with machines teaching humans

648 E.O.1 PARTICIPANTS.

We recruit 480 participants from Prolific (Palan & Schitter, 2018). We filtered the participant pool by country of residence (United States), number of completed studies (> 100), and approval rate (> 95%). Participants gave informed consent under an approved IRB protocol.

654 E.0.2 TASK.

655 We design a task for our machines to teach humans about categories, in which participants see a grid 656 of stimuli; for each cell in the grid, there is an underlying true category. Our simulated teacher model 657 selects labels based on these underlying categories, and participants see these labels with the grid. 658 Participants must then categorize all the unlabeled stimuli on the grid using the teacher's labels. We 659 do not inform participants of the number of examples per category. We investigate two structures of categories: one class per column of the grid ("cols") and one class per quadrant of the grid ("quad"). 660 Note that these categories induce labeling functions that the students *should* be able to learn; they 661 are tractable (column structure and block structure). There were two different stimuli sets. The 662 first (simple-features) is the closest analog to our simulated experiments, in which participants saw 663 a 6×6 grid with blank cells, so the features are completely expressed via the coordinates of the 664 grid. The second (*salient-dinos*) is a more rich set of stimuli, wherein participants see a 7×7 grid of 665 dinosaur ("dino") images from Malaviya et al. (2022). Dino stimuli were defined by nine different 666 features (e.g., body length, neck length, neck angle) and organized on the grid by two principle 667 components of those underlying features. For a visualization of the participant's view, see Figure 668 7. For each condition, different teachers were generated from our model, sampling across varying 669 levels of alignment. This structure leads to 24 different conditions (2 stimuli sets \times 2 category 670 structures \times 6 teacher alignment levels) for which we collect 20 participants each.

671 672

683

684 685

692

693 694

696 697

698

699

700

701

653

E.0.3 MODELS AND EXAMPLE SELECTION.

673 We employ the same model types as in our simulations. Teachers are self-centered and assume 674 that students are 1-NN classifiers¹. In both settings, we assume the representations of teachers and 675 students can be expressed through their two-dimensional grid locations. For the simple-grid setting, 676 there are no features for the human to use for their categorization beyond grid cell location; and 677 in the salient-dinos setting, features were defined by two principal components (which we can use 678 as grid coordinates). We again induce representation misalignment between teacher and student by 679 shuffling the stimuli on the grid. We sample a set of teachers spanning a range of representational alignments. We select a single set of points for each teacher, assuming the teacher has perfect 680 accuracy. We explore alternate labeling functions to simulate alternate teacher error rates post-hoc 681 (see Appendix G). 682

E.0.4 ADDITIONAL RESULTS ON MACHINES TEACHING HUMANS

We present the correlations between average human and student teacher accuracy in Table 1.

	Quadrants	Columns	Both
simple-features	0.91 (<i>p</i> =0.013)	0.59 (<i>p</i> =0.221)	0.59 (<i>p</i> =0.054)
salient-dinos	0.52 (<i>p</i> =0.286)	0.86 (<i>p</i> =0.027)	0.63 (<i>p</i> =0.037)

Table 1: Pearson correlations (with associated *p*-values) of average human student accuracy and representational alignment of the machine teacher across the various conditions.

F ADDITIONAL DETAILS ON CLASSROOM SIMULATIONS

Additional details on classroom pool construction We explore two different pools of students and teachers: (i) unstructured pools spanning a range of representational alignments and error rates, and (ii) clustered sets of students and teachers. For the latter, we construct a generative model over

¹We acknowledge such an assumption is highly simplistic for students and encourage future work to explore alternate models of students.

702 student-teacher populations wherein we have a set of clusters, with a fixed number of students per 703 cluster share similar representations. We sample one similar teacher from each cluster with some 704 error rate (sampled from a uniform distribution over 0 to 0.5). We then deliberately downsample 705 from the available teachers to simulate the case where some students are representationally distinct 706 from the other students and available pool of teachers. Additional details are included in Appendix F. For each experiment, we sample 10 different teacher pools. We additionally compute the proportion 707 of students who achieve "passing" marks (set to a moderately high threshold of 45% accurate, given 708 chance guessing is 16.6% on our 6x6 grid). We also note that we focus here on row-based labels (a 709 new f). 710

711

Utility curve over classrooms The utility curves that we construct in Section 2 and E were always 712 constructed with respect to a single student (in a respective, "dyadic"2). However, in our classroom 713 settings, we also corrupt the students' representations to simulate representational diversity. We sam-714 ple a new utility curve, wherein, for each teacher parameterization (same error rate parameterization 715 as above, with representation corruptions now in increments of 0.1), we sample 10 different student 716 corruptions ranging over 0 to 0.9 in increments of 0.1. We build this curve only for the column label 717 type. We then bucketize the teacher error rate as well as the representation alignment such that we 718 can index into the curve to extract an "expected average performance" for any student-teacher pair. 719

- Constructing classroom pools We construct two different classroom pools in Section 2.3: un structured and structured. Here, we provide additional details on how we sampled students and
 teachers for each pool type.
- 723
- **Unstructured pool** All students and teachers are sampled independently. We sample 1000 students and 30 teachers with corruption levels (pairwise swaps) sampled from a beta distribution ($\alpha = 1.5, \beta = 2.5$) to ensure that we have some students that are reasonably aligned. We sample teacher error rate uniformly over the range 0 0.5.
- 728

740

Structured pool In the structured setting, we construct *clusters* of similar students and teachers. 729 We prespecify a number of clusters M and number of students per cluster. Clusters are designed to 730 span a range of levels of representation alignment over the "original" grid. We loop over possible 731 representation alignments corruptions ranging from 0 to 1 in increments of 1/M. For each cluster, 732 we sample a "seed" student using that corruption level. We then sample students on top of this 733 cluster with a representation corruption of 0.01 on top of the base student to ensure students share 734 similar (but some variation) in their representation. For each cluster, we sample a teacher with error 735 rate uniformly from 0 - 0.5 and representation with a similar slight possible corruption (sampled 736 uniformly from 0 - 0.01) on top of the seed student, thereby ensuring that there would be a repre-737 sentationally similar teacher for each student in each cluster if provided. However, to simulate gaps in coverage of particular representation characterizations, we randomly drop some teachers from the 738 pool. 739

- 741Additional classroom matching resultsWe include additional results into classroom matching742in Tables 2 and 3 and a relationship between group size and learning outcomes in Figure 5.
- Generalization to the dino stimuli We explore generalization of our utility curve constructed in the simple-features setting to our salient-dino stimuli. We repeat our two different pool types, which we depict in Tables 4 and 5, respectively. We find that our utility curves generalize nicely to different grid sizes and stimuli type, yielding student outcomes that on average appear to boost student accuracy particularly for the students in the top-performing group than baselines which do not account for representation misalignment (MOOC).

750
751Additional details on student-centric teacherIn contrast to our self-centered teacher, our
student-centric teacher does not use its own representation to select examples to provide to the
student. Instead, the student-centric teacher is endowed with an *inner optimization loop* over the
students assigned to it, whereby the teacher loops T times over "simulated students" (which we call
the "inner loop") and randomly selects one point per category (using the teacher's believed class

²Pairing two agents – one student and one teacher.

7.50					
757	Method	Avg Acc	Bottom 10%	Top 10%	Pass Rate
758	Random	0.33 ± 0.01	0.21 ± 0.01	0.49 ± 0.01	0.12 ± 0.02
759	Min Err	0.37 ± 0.01	0.26 ± 0.01	0.53 ± 0.03	0.18 ± 0.04
760	Utility	$\textbf{0.38} \pm \textbf{0.01}$	$\textbf{0.27} \pm \textbf{0.01}$	$\textbf{0.55} \pm \textbf{0.03}$	$\textbf{0.20} \pm \textbf{0.04}$
761	Optimal	0.43 ± 0.01	0.32 ± 0.01	0.60 ± 0.02	0.32 ± 0.04

Table 2: Student learning outcomes (accuracy) from different classroom matching approaches in the structured pool setting. Higher is better for all metrics. \pm indicates standard errors computed over 40 sampled pools and associated assignments. We compute the average student performance across all N = 1000 pooled students (paired with potentially M = 30 teachers), as well as accuracy over the bottom and top 10% of students in each matching, respectively. We additionally compute the proportion of students who achieve "passing" marks (set to a moderately high threshold of 45% accurate, given chance guessing is 16.6% on our 6x6 grid). Higher is better for all metrics. \pm indicates standard errors computed over 10 sampled pools and associated assignments.

Method	Avg Acc	Bottom 10%	Top 10%	Pass Rate
Random	0.33 ± 0.00	0.17 ± 0.01	0.52 ± 0.01	0.09 ± 0.01
Min Err	$\textbf{0.39} \pm \textbf{0.01}$	$\textbf{0.25} \pm \textbf{0.00}$	0.57 ± 0.02	0.20 ± 0.02
Utility	$\textbf{0.39} \pm \textbf{0.01}$	0.24 ± 0.00	$\textbf{0.61} \pm \textbf{0.02}$	$\textbf{0.23} \pm \textbf{0.02}$
Optimal	0.49 ± 0.00	0.36 ± 0.00	0.71 ± 0.02	0.54 ± 0.01

Table 3: Student learning outcomes (accuracy) from different classroom matching approaches in the unstructured pool setting. \pm indicates standard errors computed over 40 sampled pools and associated assignments.



Figure 5: (Left:) Group sizes from greedily incorporating the lowest performing students into the classroom of a single student-centric teacher. (Right:) Average accuracy gains (out of 1.0) in performance for students grouped with the student-centered teacher, on top of what they would have achieved from a self-centered teacher. Error bars are standard errors over 20 seeds of student-centric teacher groupings for a sampled structured pool of students and teachers.

Method	Avg Acc	Bottom 10%	Top 10%	Pass Rate
Random	0.29 ± 0.00	0.16 ± 0.00	0.47 ± 0.01	0.04 ± 0.00
Min Err	0.35 ± 0.01	$\textbf{0.22} \pm \textbf{0.00}$	0.55 ± 0.04	0.13 ± 0.03
Utility	$\textbf{0.36} \pm \textbf{0.01}$	$\textbf{0.22} \pm \textbf{0.00}$	$\textbf{0.62} \pm \textbf{0.04}$	$\textbf{0.16} \pm \textbf{0.02}$
Optimal	0.44 ± 0.01	0.31 ± 0.00	0.69 ± 0.03	0.33 ± 0.01

Table 4: Student learning outcomes (accuracy) from different classroom matching approaches in the unstructured pool setting for the dino stimuli. We again compute the student performance across all N = 1000 pooled students (paired with potentially 30 teachers). Error bars are again computed over 40 different sampled pools.

810	Method	Avg Acc	Bottom 10%	Top 10%	Pass Rate
811	Wiethou	ing nee	Dottolli 1070	100 1070	1 doo rate
812	Random	0.30 ± 0.01	0.18 ± 0.01	0.45 ± 0.02	0.06 ± 0.01
813	Min Err	0.34 ± 0.01	$\textbf{0.23} \pm \textbf{0.01}$	0.51 ± 0.05	0.10 ± 0.03
814	Utility	$\textbf{0.36} \pm \textbf{0.01}$	$\textbf{0.23} \pm \textbf{0.01}$	$\textbf{0.57} \pm \textbf{0.05}$	$\textbf{0.15} \pm \textbf{0.03}$
815	Optimal	0.39 ± 0.01	0.28 ± 0.01	0.59 ± 0.04	0.18 ± 0.03

816 Table 5: Student learning outcomes (accuracy) from different classroom matching approaches in the 817 structured (clustered) pool setting for the dino stimuli. We again have 10 representationally distinct 818 clusters, each with 50 students, and sample 5 available teachers across the clusters. 819

821 - the teacher may not know the true categories) and measures the expected performance of each student if that set of examples were revealed. Note, the teacher computes the expected accuracy of 822 each student using against its belief of the true categorization (which may be incorrect). The teacher 823 then chooses the set of examples that attains the highest average accuracy over students. Here, we 824 set T to 100; exploring the impact of varied T is a sensible next step. Exploration of alternate opti-825 mization functions, e.g., optimizing over the minimum attained performance over the students in the 826 teacher's classroom rather than average classroom performance, as well as exploring different kinds 827 of simulated students (here, we assume the teacher's have the right model of each student) are also 828 ripe ground for future work. 829

We explore the effect of student-centric teachers by appending a second stage to our matching pro-830 cedure. After matching using our utility curve (as noted above), we greedily attempt to pair the 831 lowest performing students with a student-centric teacher who chooses points by optimizing for 832 the students in their pool (i.e., taking the students' representations into account). We continue in-833 corporating the next lowest-performing students into the student-centric teacher's classroom until a 834 student's attained accuracy with the original pairing is not improved by the student-centric teacher. 835 We apply our procedure to the clustered pool structure noted above and find that it is beneficial to 836 continue adding students up to a point: if the teacher is an expert (zero error rate), we can add all stu-837 dents from one cluster before we see detrimental performance across the pool of students assigned 838 to said teacher. As the student-centric teacher's error rate increases, fewer students can be pooled 839 before performance dropoff (see Appendix Figure 5).

840 These results indicate the student-centric teachers can cover students who are representationally dis-841 tinct and help boost their learning outcomes. However, classroom size matters, corroborating prior 842 works in machine and human teaching (Frank, 2014; Yeo et al., 2019; Ma et al., 2018; Zhu et al., 843 2018). In the next section of the Appendix, we conduct a deeper dive into the relationship between 844 classroom size and student outcomes in our setting when student-centric teachers are available. Herein, we see that teachers who may try to overalign to all students at once in a large classroom 845 induce poorer outcomes for the classroom writ large. 846

847 848

849

851

855

856

857 858

860

G ADDITIONAL HUMAN EXPERIMENT DETAILS

850 **Participant recruitment and compensation** Participants were recruited from Prolific and were paid \$12/hr plus a 10% bonus if they responded reasonably (i.e., did not select labels randomly or 852 choose the same label for all stimuli). The research did not contain risks to participants, and they 853 were able to opt out at any time. The institution of the principal investigator obtained IRB approval 854 for this experiment, and participants gave informed consent under this protocol.

Task instructions We include the full set of instructions provided to participants in Figure 6 and sample interfaces in Figure 7.

Further analyses 859

861 **Simulating teacher error in human experiments** All human experiments were run with machine teachers set to zero error, as collecting all combinations of teacher error and representational align-862 ment would be prohibitively expensive. Instead, we simulate the effect of teacher error in a post-hoc 863 analysis by corrupting the true underlying labels in the same way we corrupted the teacher labels for

864	
865	
866	
867	You are a Student in our Teacher-Student interaction experiment. You will be paired with a
868	Teacher.
869	You will both be shown images that represent stick figure dinosaurs on a 7-by-7 grid.
870	The grid is split into <mark>7 categories</mark> of dinosaurs.
871	Every dinosaur on the grid is in one of these categories (from A to G).
872	
873	Your goal is to guess the category of every dinosaur on the grid.
874	The Teacher's goal is to help you guess the categories of dipercurs correctly by revealing one
875	The reacher's goal is to help you guess the categories of dinosaurs correctly by revealing one
876	label for each type to you.
877	
878	You will receive <mark>bonus compensation based on how many labels you guess correctly,</mark> so
879	please do your best.
880	
881	Next
882	

885 Figure 6: Experiment instructions displayed to all participants, introduced paragraph by paragraph. 886 The only changes to instructions were to modify the type of stimuli ("empty cells", "images that represent stick figure dinosaurs"), size of the grid $(6 \times 6, 7 \times 7)$, the number and names of categories 888 (4; A-D or 6/7; A to F/G).



912 913

883 884

887

Figure 7: Above are two example views of the experiment. All participants, after viewing the 914 instructions in Figure 6 were taken to a page that contained a grid and the labeled stimuli. They 915 were asked to categorize stimuli via a dropdown menu selection. Finally, they rated their confidence 916 using a scale below the stimulus grid. Left: salient-dinos, 7 ("col") categories, medium-alignment 917 teacher. Right: simple-features, 4 ("quad") categories, high-alignment teacher.



Figure 8: Average human student classification accuracy at various levels of representational alignment. Error bars correspond to one standard error. (Left:) Results from simple-features setting.
(Right:) Results from salient-dinos setting. (Top:) One class per quadrant. (Middle:) One class per column (6 for simple-features, 7 for salient-dinos). (Bottom:) Combined results.

the simulation experiments (i.e., error rate corresponds to the probability with which we flip each true label to be a different label). Human student accuracy was then recomputed against these corrupted true labels. The original human student results with no simulated teacher error are reported in Figure 8.