

# Towards Interactive Language Modeling

Anonymous ACL submission

## Abstract

Interaction between caregivers and children plays a critical role in human language acquisition and development. Given this observation, it is remarkable that explicit interaction plays little to no role in artificial language modeling—which also targets the acquisition of human language, yet by artificial models. Moreover, an interactive approach to language modeling has the potential to make language models substantially more versatile and to considerably impact downstream applications. Motivated by these considerations, we pioneer the space of interactive language modeling. First we present a road map in which we detail the steps that need to be taken towards interactive language modeling. We then lead by example and take the first steps on this road map, showing the initial feasibility of our approach. As such, this work aims to be the start of a larger research agenda on interactive language modeling.

## 1 Introduction

Interaction between children and more advanced language interlocutors (such as caregivers) plays an important role in many theories and studies on human language acquisition (e.g., Bruner, 1985; Clark, 2018). For example, although culturally dependent (Shneidman and Goldin-Meadow, 2012) and with the precise effects still up for discussion (Cristia et al., 2019), caregivers can communicate with their children in Child Directed Speech. In turn, children can for example experiment with the meaning of words, to elicit a response from their caregivers (Gillis and Schaeerlaekens, 2000).

Despite the importance of interaction in human language acquisition, interaction plays little to no role in artificial language modeling. This is remarkable, as language modeling also has the objective to learn human language, albeit with artificial models. Instead, current state-of-the-art language models (LMs) take large amounts of text as input, and are tasked to predict the next or masked words (e.g.,

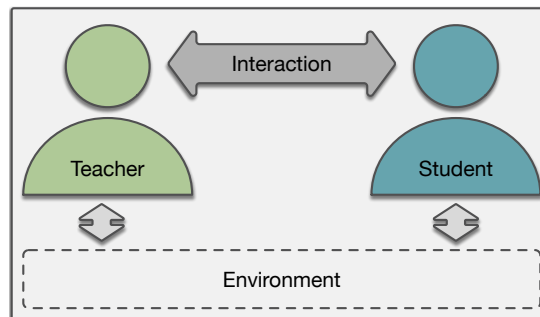


Figure 1: Teacher-Student setup for interactive language modeling.

Devlin et al., 2018; Brown et al., 2020). The learning signal only comes from a cross-entropy loss that indicates whether a prediction was correct. Although this setup has shown to be effective, from the perspective of human language acquisition it appears very unnatural. This gives rise to the motivation to investigate other, more natural approaches to language modeling, such as the interactive perspective that we propose in this paper.

Specifically, we structure our proposal according to a teacher-student setup. Figure 1 depicts a high level overview. In this setup we distinguish four main parts: *the teacher*, whose role is inspired by the caregiver in the human language acquisition, *the student*, who resembles the child, *the interaction* between the teacher and the student and *the environment* that they both share (such as the language that needs to be learned by the student). We motivate and detail our setup further in Section 3.

An interactive approach to language modeling is not only interesting from the perspective of human language acquisition. Explicitly allowing for interaction also has the potential to make language modeling more efficient and versatile. For example, a teacher can adapt its input to a student based on the specific feedback signals it receives from the student, and a teacher that is fluent in one domain can teach the specifics of that domain to a student trained on another domain, and vice versa. More-

071	over, an interactive approach to language modeling	<b>2.3 Active Learning</b>	118
072	has the potential to impact downstream applica-	Active Learning (AL) (Cohn et al., 1996) is an ap-	119
073	tions, for example for foreign language teaching	proach in which a learner (the model to be trained)	120
074	apps where a student can be replaced by a human.	actively selects which data it can most effectively	121
075	In this paper we pioneer the space of interactive	be trained on. That is, where CL is often more	122
076	language modeling. Specifically, we contribute:	associated with choosing a teaching strategy, AL	123
077	<b>C1</b> We define the objective of interactive language	is rather focused on the student side. AL is of-	124
078	modeling;	ten used to efficiently label data in a low resource	125
079	<b>C2</b> We present a road map that details the steps	setting (e.g., Reichart et al., 2008; Dor et al., 2020).	126
080	that need to be taken towards this objective;		
081	<b>C3</b> We take the first steps on this road map, which	<b>2.4 Continual Learning</b>	127
082	show the initial feasibility of our approach.	In Continual Learning, or life-long learning, the	128
083	By doing so we aim to start a larger research	aim is to train a model in an online fashion, i.e., on	129
084	agenda on interactive language modeling.	a continuous stream of data, whilst avoiding <i>catas-</i>	130
085	<b>2 Related Work</b>	<i>trophic forgetting</i> (McCloskey and Cohen, 1989;	131
086	Over the years many different types of learning	French, 1999). This makes models versatile to	132
087	strategies have been proposed for artificial model-	an ever changing world. Some recent work has	133
088	ing. Below we describe a number of them that are	focused on types of Continual Learning for large	134
089	particularly related to the current work.	LMs (e.g., Lazaridou et al., 2021; Jin et al., 2021).	135
090		We envision interactive language modeling to play	136
091	<b>2.1 Interactive Language Learning in NLP</b>	an important role in life-long learning in the future.	137
092	Recently, a number of studies have focused on		
093	interactive language learning. Stein et al. (2021)	<b>3 A Road Map towards Interactive</b>	138
094	learn logical semantic representations in an inter-	<b>Language Modeling</b>	139
095	active way. Nikolaus and Fourtassi (2021) propose	In this section we present a general road map to-	140
096	a proof of concept to model perception and pro-	wards interactive language modeling.	141
097	duction based learning of semantic knowledge ac-	<i>Our objective is to build an automated teacher-</i>	142
098	quisition in children. Kiseleva et al. (2021) take	<i>student loop for language modeling that attains</i>	143
099	an interactive approach to language <i>understand-</i>	<i>good performance in the student for a fixed (low)</i>	144
100	<i>ing</i> in a recent NeurIPS challenge. To the best of	<i>number of bits transmitted in the interactions.</i>	145
101	our knowledge, none of the existing works have	We propose a teacher-student loop as this format	146
102	focused specifically on language modeling.	closely resembles caregiver-child interactions. In	147
103		Section 1 and Figure 1 we already introduced a	148
104	<b>2.2 Curriculum Learning</b>	high level overview of this setup and its four main	149
105	Curriculum Learning (CL) (Bengio et al., 2009) is	components: (1) <i>the teacher</i> , (2) <i>the student</i> , (3) <i>the</i>	150
106	an approach to learning in which data samples are	<i>interaction</i> and (4) <i>the environment</i> . Generally, in	151
107	presented in a meaningful order—typically in order	this setup teachers transmit language data to their	152
108	of complexity—motivated by the idea that humans	students, according to a certain budget (“a (low)	153
109	learn in a similar way. Bengio et al. show the	fixed number of bits”). Having this budget forces	154
110	effectiveness of CL on a number of tasks, among	the teacher to actively choose a learning strategy, as	155
111	which a classical approach to language modeling.	just sending all data that is available to the teacher	156
112	More recently, a number of studies have shown	would not be allowed. Students have the objective	157
113	the effectiveness of CL for (fine-tuning) LMs (Xu	to learn the language and they send a signal back	158
114	et al., 2020; Zhang et al., 2021), although other	that informs their teacher of their performance, e.g.,	159
115	studies have shown that not all intuitive curricula	a score on an exam. This interaction takes place in	160
116	are also effective (Liu et al., 2019). Matiisen et al.	an environment, e.g., a common language.	161
117	(2019) propose a teacher-student framework for	In Table 1 we present the road map that we envi-	162
	automatic CL for the addition of decimal numbers	sion towards interactive language modeling. This	163
	and navigation in Minecraft.	road map works as follows. For each of the four	164
		mentioned components we detail steps that	165

need to be taken. We also add a fifth component: the evaluation of the setup. Each component has different aspects (bold-faced in Table 1). For example, for the *teacher* we can focus on how it can access the data that it can transmit to the student, which we call “ways of speaking” in Table 1. Another aspect of the teacher side focuses on what we call the “degree of awareness”, which entails different ways in which the teacher can remember different aspects of the teaching loop. In a similar fashion we fill in the remaining components in the table. We focus on text as a single modality and acknowledge grounded interactive language modeling as an interesting future research direction.

On our road map there are multiple ways to reach the destination. For example, one can focus on taking a few steps for each of the components, or to take many steps for only one or a few of the components. Moreover, although mostly structured in increasing degree of complexity, this does not always hold for all individual steps in the table. For example, zooming in on the “degrees of awareness” for the teacher again, one could imagine an example where a teacher does not have an explicit memory buffer of what it sent to the student before, but does have an explicit way of remembering what the student’s fine-grained capabilities are, as well as the other way around.

In the remainder of this work we take the first steps on the road map. We focus on the teacher side, i.e., learning the correct didactic approach.

## 4 Taking the First Steps on the Road Map

Figure 2 shows how we adapt the general setup from Figure 1 to take the first steps on the road map. Here we describe each modification per component: *the teacher*, *the student*, *the interaction*, *the environment* and *the exam* that the student takes.

### 4.1 The Teacher

In this work we focus on the teacher side. The role of the teacher is to transmit language data that will optimally help the student to learn the language. Figure 2 shows that we train the teacher to do this in a number of time steps. At each of these steps a teacher samples data from a larger language data set according to a fixed budget. We discuss the specifics of the sampling function below. To reduce the variance in the teacher’s learning process we repeat this process for multiple students, i.e., a teacher selects  $N$  “lessons” for  $N$  students. Due

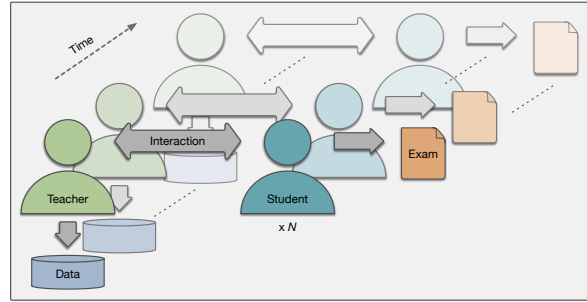


Figure 2: Teacher-student loop as used in this work.

to the stochasticity of the sampling process, each student has the potential to be trained on a slightly different part of the data. Because we use a multi-processing setup we can train multiple students on a single GPU. Hence, using multiple students does not drastically increase the computational cost.

#### 4.1.1 Knowing the Language

The teacher is modeled as a native speaker of the language that it needs to teach. We represent the teacher’s language understanding with a pretrained causal Transformer LM (Vaswani et al., 2017). We pretrain this model on a *different* subset of the data than the teacher can select from for the students, and thus we ensure that we measure whether a teacher can teach a language as a whole, and not only a particular subset that it was trained on itself.

#### 4.1.2 Selecting the Data

We use REINFORCE (Williams, 1992) with entropy regularization (Mnih et al., 2016) to learn the teacher’s didactic approach.<sup>1</sup> We want to optimize the teacher’s policy such that it learns to select the optimal data to train the student on, given a predefined budget. The policy is a one-layer feed forward neural network, that outputs a score for each sentence, i.e., the teacher’s policy network takes a sentence embedding as input, based on the pretrained Transformer LM that we use to represent the teacher’s language understanding. An action is modeled as selecting  $k$  sentences from the larger data set, where  $k$  is a predefined teacher budget. We use the GumbelTopK trick (Vieira, 2014; Kool et al., 2019) to sample  $k$  sentences without replacement, based on the teacher policy’s output scores. We compute the log probabilities (needed to compute the loss) for each sample by adding the log

<sup>1</sup>We also experimented with gradient-free optimization approaches such as the ones implemented in Nevergrad (Rapin and Teytaud, 2018), but found REINFORCE to be more flexible in our case and therefore a better fit for our needs.

Teacher	Student
<p><b>Ways of speaking</b></p> <ul style="list-style-type: none"> <li>• Select data from bin;</li> <li>• Generate data with own language model.</li> </ul> <p><b>Degrees of awareness</b></p> <ul style="list-style-type: none"> <li>• (No) memory buffer of what has been sent to the student and being able to act on it (see <i>Interaction</i> cell);</li> <li>• (No) explicit way of remembering what the student’s fine-grained capabilities are and being able to act on it (see <i>Interaction</i> cell).</li> </ul>	<p><b>Ways of speaking</b></p> <ul style="list-style-type: none"> <li>• Generate language data in a standard LM fashion;</li> <li>• Actively experiment with language generation to elicit direct feedback from the teacher (see also <i>Interaction</i> cell).</li> </ul> <p><b>Degrees of using the teacher data</b></p> <ul style="list-style-type: none"> <li>• Use all data received from the teacher;</li> <li>• Actively select data that is useful;</li> <li>• Actively know when to stop training (for example to avoid overfitting).</li> </ul>
Interaction	Environment
<p><b>Teacher side</b></p> <ul style="list-style-type: none"> <li>• Send all data at once;</li> <li>• Send data in batches, based on student feedback (see below). Batches can be as small as single utterances, after which the student sends an utterance back, like in real human-to-human interaction (see below);</li> <li>• Send (mid-term) exams.</li> </ul> <p><b>Student side</b></p> <ul style="list-style-type: none"> <li>• Send a single average exam score back to the teacher;</li> <li>• Send a fine-grained exam score back, e.g., <ul style="list-style-type: none"> <li>– score per item on the exam set;</li> <li>– (average) scores of different components (tasks) of the exam(s)</li> </ul> </li> <li>• Ask for feedback, for example by actively experiment with language generation for the teacher to judge (‘generate own exam’).</li> </ul>	<p><b>Language</b></p> <ul style="list-style-type: none"> <li>• Artificial languages, in increasing level of difficulty in terms of complexity, e.g., <ul style="list-style-type: none"> <li>– random language;</li> <li>– different types of structures;</li> <li>– different vocabulary sizes;</li> </ul> </li> <li>• Subset of human language, e.g., in terms of <ul style="list-style-type: none"> <li>– semantics (e.g., different domains)</li> <li>– syntax (e.g., different grammatical structures)</li> <li>– pragmatics</li> </ul> </li> <li>• Unrestricted human language.</li> </ul> <p><b>Task</b></p> <ul style="list-style-type: none"> <li>• <i>Teacher</i>: Learn to select or generate the optimal data such that the student performs well on the exam set (see cell below);</li> <li>• <i>Teacher</i>: Learn to adapt to different types of students, e.g., <ul style="list-style-type: none"> <li>– architectural differences</li> <li>– different prior knowledge (be aware of catastrophic forgetting in neural networks)</li> </ul> </li> <li>• <i>Student</i>: Learn to adapt to different types of teachers (didactic strategies).</li> </ul>
Evaluation / Exam	
<p><b>Teacher</b></p> <ul style="list-style-type: none"> <li>• Accuracy in selecting the optimal teaching protocol</li> </ul> <p><b>Student (Exam / Feedback for teacher)</b></p> <ul style="list-style-type: none"> <li>• General performance, measured in perplexity;</li> <li>• Performance on specific tasks, such as <ul style="list-style-type: none"> <li>– Subset of the data known to the teacher (e.g., specific domain or (grammatical) structure)</li> <li>– BLIMP (Warstadt et al., 2020);</li> <li>– BIG-Bench (<a href="https://github.com/google/BIG-bench">https://github.com/google/BIG-bench</a>).</li> </ul> </li> <li>• Scores either as an average of more fine-grained (see <i>Interaction</i> cell).</li> </ul>	

Table 1: Road map to interactive language modeling.

250	probabilities of each element in the sample. We	$V_2 = \{k, l, m, n, o, p, q, r, s, t\}$ . We construct sen-	296
251	explain the rationale behind this in Appendix A.	ences by randomly sampling from these sets. Sen-	297
252	<b>4.2 The Student</b>	tences consist either of tokens only from $V_1$ or	298
253	As the teacher is the main focus of our work, we	of tokens only from $V_2$ . Sentences have an equal	299
254	choose to keep the student side simple. We repre-	length of 10 tokens each. Half of the data set that	300
255	sent the student as a causal Transformer LM, that	the teacher can choose from consists of $V_1$ sen-	301
256	we train on the data that it receives from the teacher.	tences, the other half consists of $V_2$ sentences. The	302
257	<b>4.3 The Interaction</b>	teacher’s LM is trained on a similarly constructed	303
258	Following Table 1, the teacher sends all selected	data set, yet consisting of different sentences. The	304
259	data to the student at once. The student uses this	student’s exam set consists of sentences from only	305
260	data to train its LM and takes an exam after a prede-	one of the vocabularies, $V_1$ in our case. These are	306
261	defined number of updates. The average exam score	different sentences than in the training set, i.e., the	307
262	is sent back to the teacher as feedback. We use	teacher cannot simply sample the exam set to train	308
263	the student’s last model checkpoint to compute the	the student. Hence, the optimal teaching strategy	309
264	scores (as opposed to the best checkpoint on a vali-	is to present the student with sentences from the	310
265	dation set), to ensure that the learning signal for the	exam vocabulary. We confirm this in our baseline	311
266	teacher is restricted to the student’s performance	experiments that we present in Section 5.4.	312
267	on the exam set, i.e., we do not expect teachers to	<b>5.2 Task 2 – Teaching Different Structures</b>	313
268	reverse the learning process of the students (just	For this task we do not use different vocabularies,	314
269	like caregivers cannot do this for their children).	but different sentence structures. All our sentences	315
270	<b>4.4 The Environment</b>	are constructed with $V_1$ and are between 2 and 10	316
271	Following Table 1, we design a number of artifi-	tokens long. We use two different structures: single	317
272	cial languages to test our approach on (see Sec-	repetitions and double repetitions. In the case of	318
273	tion 5 for details). Using artificial languages is a	the single repetitions two identical tokens never	319
274	well-tested approach to study the behavior of neu-	occur next to each other, whereas in the case of	320
275	ral networks (e.g., Batali, 1994; Wiles and Elman,	double repetitions tokens are sampled in pairs:	321
276	1995; Rodriguez et al., 1999; Gers and Schmidhu-	<i>Structure 1</i> - Single repetitions: $(xy)^n$	322
277	ber, 2001; Rodriguez, 2001; Hupkes et al., 2018;	<i>Structure 2</i> - Double repetitions: $(xx)$ or $(xxyy)^n$	323
278	Lake and Baroni, 2018; Saxton et al., 2019; Hupkes	The data set that the teacher can sample from con-	324
279	et al., 2020; Rodríguez Luna et al., 2020; van der	sists for 20% of sentences with Structure 1 and for	325
280	Wal et al., 2020; Chaabouni et al., 2021; Dagan	80% of sentences of Structure 2. The exam set	326
281	et al., 2021). Using artificial languages gives us the	consists of sentences with Structure 1. We opt for	327
282	control we need to design our experiments in such	this way of splitting the data, as we found that a	328
283	a way that we can correctly interpret the results.	student performs quite well when trained on data	329
284	<b>4.5 The Exam</b>	consisting half of Structure 1 and half of Structure	330
285	The exam is a held-out set over which we compute	2. Having an unequal split thus allows us to make	331
286	the student’s perplexity. The details of the exam	sure that we can appropriately distinguish a learned	332
287	are task dependent and we discuss these next.	didactic approach from a random one. For this task	333
288	<b>5 Experiments</b>	the optimal teaching strategy is to select sentences	334
289	We test our proposed setup on a number of settings	with the exam structure, as we confirm with our	335
290	and tasks, that we describe in this section.	baseline experiments that we present in Section 5.4.	336
291	<b>5.1 Task 1 – Teaching Different Domains</b>	<b>5.3 Training Details</b>	337
292	For this task we design a language consisting of	The teacher LM is trained on 100 unique sentences	338
293	two strictly separated vocabularies, loosely repre-	till convergence. The dataset the teacher can sam-	339
294	sending two different domains in natural language.	ple from for the student consists of 100 different	340
295	Specifically, $V_1 = \{a, b, c, d, e, f, g, h, i, j\}$ , and	unique sentences. The exam consists of 10 unique	341
		sentences and we set the teacher budget to 10 as	342
		well. We run our experiments with five different	343
		random seeds and report the averages and standard	344
		deviations. We use the negative perplexity of the	345

student on the exam as reward for the teacher. We experiment with two different sentence embeddings for the teacher: average word embeddings and the average of the last hidden layer. We train students for a predefined number of steps that we determine by inspecting the loss and perplexity curves of training an LM once before the actual experiments. We base the threshold on when a student LM starts to overfit, so that a teacher can get clear feedback signals. We set this value to 400 for Task 1 and 300 for Task 2. Automatically determining when the students stops training is an important avenue for future work (Table 1). We use Fairseq’s (Ott et al., 2019) `transformer_lm`<sup>2</sup> for the implementation of the Transformer LMs. We use up to four GPUs with 32 GB RAM per experiment. The exact number depends on the number of students per teacher, as we can fit up to 6 students on a single GPU due to our multiprocessing implementation.

## 5.4 Baseline experiments

We run three baseline experiments with three different didactic strategies: an *oracle*, *random*, and *worst case* strategy. We run the baselines for five different random seeds. In each experiment, we randomly select data according to the teacher budget. We do this five times and each time train a student LM with the selected data. The difference between baselines is the type of data that can be selected. For the oracle baseline we only select sentences that consist of the exam vocabulary (Task 1) or structure (Task 2). For the random baseline we randomly select sentences. For the worst case baseline all sentences that we select are from a different vocabulary or structure than the exam sentences.

## 6 Results

### 6.1 Task 1 – Different Domains

#### 6.1.1 Baseline Results

In Table 2 we present the results for the baseline experiments for Task 1. We report the averages and standard deviations of the perplexity on the exam set and the fraction of training sentences that consisted of the exam vocabulary. For space reasons, we report the results for two seeds per baseline: the seed with the best average perplexity and the worst. The results for all five seeds are given in Appendix B. There we also present scores for the  $n$ -gram overlap between the selected training set

<sup>2</sup>[https://fairseq.readthedocs.io/en/latest/command\\_line\\_tools.html](https://fairseq.readthedocs.io/en/latest/command_line_tools.html)

Type	Seed	Avg Perplexity	Avg train from test
<i>Rand.</i>	B	$160.9 \pm 217.7$	$0.54 \pm 0.16$
	W	$742.5 \pm 159.8$	$0.50 \pm 0.17$
<i>Orac.</i>	B	$14.99 \pm 5.364$	$1.00 \pm 0.00$
	W	$68.95 \pm 87.49$	$1.00 \pm 0.00$
<i>Worst case</i>	B	$4.78e4 \pm 2.67e4$	$0.00 \pm 0.00$
	W	$8.46e4 \pm 4.69e4$	$0.00 \pm 0.00$

Table 2: Baseline results Task 1. Averages and standard deviations reported based on five runs per seed. *Rand* is Random, *Orac* is Oracle, *B* is Best and *W* is Worst.

and the exam set. The results are as expected. The oracle baseline gives the best results, followed by the random and worst case baseline respectively.

### 6.2 Results of Training the Teacher

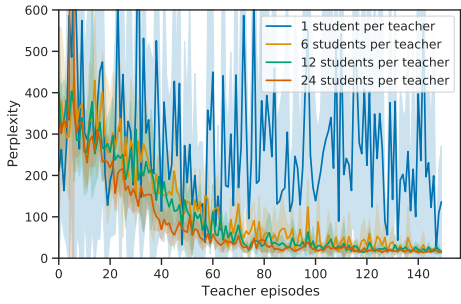
In Figure 3 we present the results for Task 1 for different numbers of students per teacher.<sup>3</sup> The teacher’s didactic strategy correctly converges to the oracle baseline. There is a clear difference between different sentence embeddings (Section 4.1.1). Both embedding types are converging, but the average hidden layer embeddings are clearly superior. We investigate this further by plotting the t-SNE embeddings (Van der Maaten and Hinton, 2008) of the different sentence embeddings in Figure 4. To prepare for Task 2, we also plot the embeddings of Task 2. The hidden layer sentence embeddings result in the clearest separation between sentences from different vocabularies or structures. Especially for Task 2, where we use the same vocabulary, this is unsurprising. From now on we opt for these sentence embeddings. Based on the results for Task 1 we opt for 12 students per teacher as a good trade-off between computational cost and convergence stability for Task 2.

### 6.3 Task 2 – Different Structures

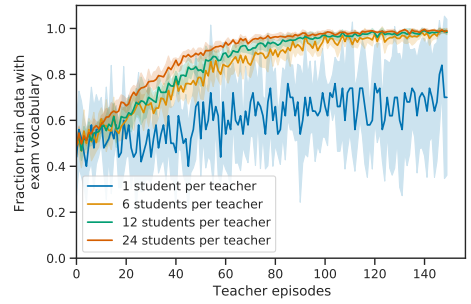
#### 6.3.1 Baseline Results

We present the results for the baseline results for Task 2 in Table 3. Again we report the results for the best and the worst seed. Full results are available in Appendix C. Similarly to the results for Task 1, we confirm that the oracle baseline performs strongest, followed by the random and worst case baseline respectively.

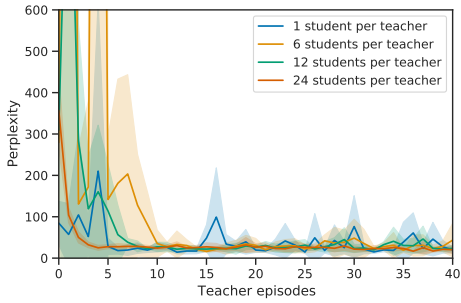
<sup>3</sup>We present plots for the  $n$ -gram overlap in Appendix D.



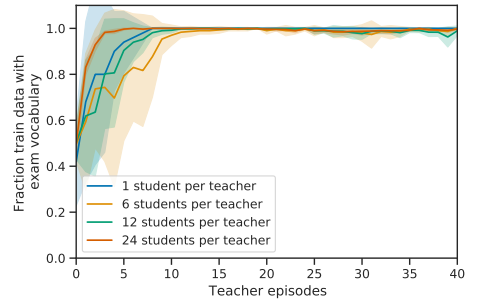
(a) Perplexity of the student on the exam data over different episodes. Average word embedding as input to the teacher's policy.



(b) Fraction training data with the exam vocabulary over different episodes. Average word embedding as input to the teacher's policy.

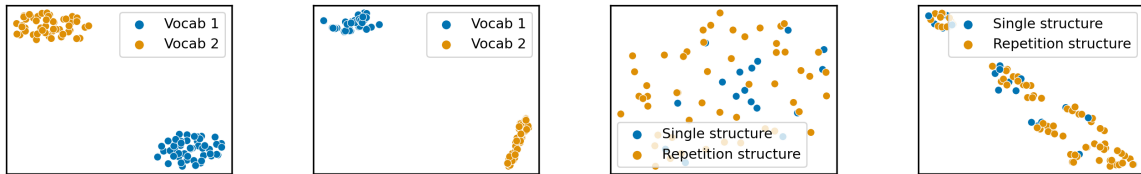


(c) Perplexity of the student on the exam data over different episodes. Average last hidden layer as input to the teacher's policy.



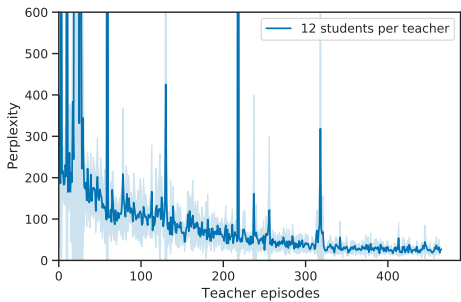
(d) Fraction training data with the exam vocabulary over different episodes. Average last hidden layer as input to the teacher's policy.

Figure 3: Results Task 1 – Different domains. Plots for different numbers of students per teacher. Results per setting reported as average and standard deviation over five random seeds. x-axis of lower plots bound to 40 as the teacher had already converged by then.

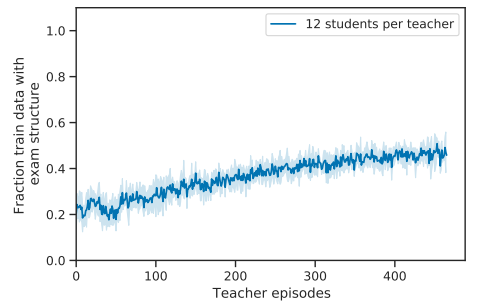


(a) Task 1 - Different vocabularies. Sentence embedding is average word embeddings. (b) Task 1 - Different vocabularies. Sentence embedding is average last hidden layer. (c) Task 2 - Different structures. Sentence embedding is average word embeddings. (d) Task 2 - Different structures. Sentence embedding is average last hidden layer.

Figure 4: T-SNE plots for different sentence representations for different tasks.



(a) Perplexity of the student on the exam data over different episodes.



(b) Fraction training data with the exam structure over different episodes.

Figure 5: Results Task 2 – Plots for 12 students per teacher. Results per setting reported as average and standard deviation over five random seeds.

Type	Seed	Avg Perplexity	Avg train from test
<i>Rand.</i>	B	119.0 ± 56.48	0.18 ± 0.04
	W	342.1 ± 241.4	0.12 ± 0.08
<i>Orac.</i>	B	6.821 ± 0.619	1.00 ± 0.00
	W	9.431 ± 3.057	1.00 ± 0.00
<i>Worst Case</i>	B	299.6 ± 124.2	0.00 ± 0.00
	W	595.3 ± 297.9	0.00 ± 0.00

Table 3: Baseline results Task 2. Averages and standard deviations reported based on five runs per seed. *Rand* is Random, *Orac* is Oracle, *B* is Best and *W* is Worst.

### 6.3.2 Results of Training the Teacher

In Figure 5 we present the results for Task 2.<sup>4</sup> Again we see that the teacher learns to gradually converge to the oracle teaching strategy, although convergence is less fast than for Task 1; we do not achieve full convergence in the number of training episodes that we run these experiments for. We postulate that this can be explained by the differences we found in Figure 4. The differences in sentence embeddings between the two different structures are clearly less apparent than between the sentences from two vocabularies. This indicates the importance of good sentence embeddings in future work. Moreover, as stated in Section 6.3, we found that transmitting roughly 50% of Structure 1 and 50% of Structure 2 also already leads to good performance. Therefore, the teacher likely needs to learn from a less distinct learning signal than in Task 1.

## 7 Implications and Outlook

We successfully took the first steps on our proposed road map. Here we want to share our learnings and the limitations of the current setup to help future research to take the next steps on the road map.

**The importance of designing experiments with interpretable outcomes.** We designed our experiments such that we knew the teacher’s oracle strategy, which allowed us to properly test our setup. However, in designing our experiments we found that finding such settings is non-trivial. For example, in a task that contains a language with multiple structures, a student might unexpectedly learn information from structure 1 that also proves useful for structure 2. This might be acceptable if one’s only objective is to obtain a good performance. How-

<sup>4</sup>We present plots for the  $n$ -gram overlap in Appendix E.

ever, in our case it is critical to be able to know that a teacher is “right for the right reasons”, which motivated our choices for the tasks and languages. **The teacher’s budget.** Following our objective, we designed our experiments in such a way that the teacher was given a budget that limits the amount of data it can send to the student. As mentioned in Section 5.3, we confirmed that the student’s learning converges with this budget. In follow up work we plan to investigate the importance of different budgets in more detail. One interesting direction is to give the teacher a flexible budget, i.e., such that a teacher could decide to stop training if it deems it no longer necessary for the student.

**Computational complexity.** Apart from the multiprocessing setup that allows us to train multiple students on a single GPU, we did not yet focus on the computational complexity of our approach. In the current setup many student language models need to be trained for a single teacher. In our case we deem this justifiable as we are just at the start of the road map. Moreover, once a teacher model is trained, it can be used for many different purposes. However, in future work we hope to focus on decreasing the computational complexity of our approach. One promising avenue to do this is by optimizing the learning process of the student.

## 8 Conclusion

In this paper we pioneered the space of interactive language modeling, motivated by the observation that current state-of-the-art LMs are trained in a very unnatural way, from the perspective of human language acquisition. Moreover, an interactive approach has the potential to make LMs more efficient and adaptable. Specifically, we proposed a teacher-student loop, in which the teacher is inspired by the caregiver and the student resembles the child in the human language acquisition. We presented a road map that details the steps towards interactive language modeling for each of the components of the teacher-student loop. We led by example and took the first steps on this road map, leading to a tangible proof of concept of our proposal. As such, we structured the space of interactive language modeling and aim to inspire a larger research agenda on interactive language modeling.

## 9 Ethical Impact Statement

At this point we use artificial language data only, for which we do not see any direct negative impli-



510	cations. As we move towards using real data sets,	Gautier Dagan, Dieuwke Hupkes, and Elia Bruni. 2021.	559
511	it is necessary to be aware of potential biases with	<a href="#">Co-evolution of language and agents in referential</a>	560
512	these data sets. One needs to ensure that the data is	<a href="#">games</a> . In <i>Proceedings of the 16th Conference of the</i>	561
513	not biased towards any (protected) group to avoid	<i>European Chapter of the Association for Computa-</i>	562
514	any harm. Currently, much of the NLP research	<i>tional Linguistics: Main Volume</i> , pages 2993–3004,	563
515	focuses on English as its language of interest. Our	Online. Association for Computational Linguistics.	564
516	approach is not bound to any language in particular	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	565
517	and can even be used to improve language learn-	Kristina Toutanova. 2018. Bert: Pre-training of deep	566
518	ing in a low resource setting. Once the models	bidirectional transformers for language understand-	567
519	achieve human like performance and are used for	ing. <i>arXiv preprint arXiv:1810.04805</i> .	568
520	downstream tasks and applications it is necessary	Yue Dong, Yikang Shen, Eric Crawford, Herke van	569
521	to explicitly state that language is produced by an	Hoof, and Jackie Chi Kit Cheung. 2018. Bandit-	570
522	artificial language model. However, as with all lan-	sum: Extractive summarization as a contextual ban-	571
523	guage models, misuse can still happen and it is our	dit. <i>arXiv preprint arXiv:1809.09672</i> .	572
524	responsibility as a research community, amongst	Liat Ein Dor, Alon Halfon, Ariel Gera, Eyal Shnarch,	573
525	others, to spend effort on making users aware of	Lena Dankin, Leshem Choshen, Marina Danilevsky,	574
526	these possibilities.	Ranit Aharonov, Yoav Katz, and Noam Slonim. 2020.	575
		Active learning for bert: An empirical study. In	576
		<i>Proceedings of the 2020 Conference on Empirical</i>	577
		<i>Methods in Natural Language Processing (EMNLP)</i> ,	578
		pages 7949–7962.	579
527	<b>References</b>	Robert M French. 1999. Catastrophic forgetting in con-	580
528	John Batali. 1994. Artificial evolution of syntactic ap-	nectionist networks. <i>Trends in cognitive sciences</i> ,	581
529	ptitude. In <i>Proceedings from the Sixteenth Annual</i>	3(4):128–135.	582
530	<i>Conference of the Cognitive Science Society</i> , pages	Felix A Gers and Jürgen Schmidhuber. 2001. LSTM	583
531	27–32. Lawrence Erlbaum Associates Hillsdale, NJ.	recurrent networks learn simple context-free and	584
		context-sensitive languages. <i>Neural Networks, IEEE</i>	585
532	Yoshua Bengio, Jérôme Louradour, Ronan Collobert,	<i>Transactions on</i> , 12(6):1333–1340.	586
533	and Jason Weston. 2009. Curriculum learning. In	Steven Gillis and Annemarie Schaerlaekens. 2000.	587
534	<i>Proceedings of the 26th annual international confer-</i>	<i>Kindertaalverwerving: Een handboek voor het Ned-</i>	588
535	<i>ence on machine learning</i> , pages 41–48.	<i>erlands</i> .	589
		Dieuwke Hupkes, Verna Dankers, Mathijs Mul, and Elia	590
536	Tom B Brown, Benjamin Mann, Nick Ryder, Melanie	Bruni. 2020. Compositionality decomposed: how do	591
537	Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind	neural networks generalise? <i>Journal of Artificial</i>	592
538	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	<i>Intelligence Research</i> , 67:757–795.	593
539	Askell, et al. 2020. Language models are few-shot	Dieuwke Hupkes, Sara Veldhoen, and Willem Zuidema.	594
540	learners. <i>arXiv preprint arXiv:2005.14165</i> .	2018. Visualisation and ‘diagnostic classifiers’ reveal	595
		how recurrent and recursive neural networks process	596
541	Jerome Bruner. 1985. Child’s talk: Learning to use	hierarchical structure. <i>Journal of Artificial Intelli-</i>	597
542	language. <i>Child Language Teaching and Therapy</i> ,	<i>gence Research</i> , 61:907–926.	598
543	1(1):111–114.	Xisen Jin, Dejiao Zhang, Henghui Zhu, Wei Xiao,	599
		Shang-Wen Li, Xiaokai Wei, Andrew Arnold, and	600
544	Rahma Chaabouni, Roberto Dessì, and Eugene	Xiang Ren. 2021. Lifelong pretraining: Continu-	601
545	Kharitonov. 2021. Can transformers jump around	ally adapting language models to emerging corpora.	602
546	right in natural language? assessing performance	<i>arXiv preprint arXiv:2110.08534</i> .	603
547	transfer from scan. <i>arXiv preprint arXiv:2107.01366</i> .	Julia Kiseleva, Ziming Li, Mohammad Aliannejadi,	604
		Shrestha Mohanty, Maartje ter Hoeve, Mikhail Burt-	605
548	Eve V Clark. 2018. Conversation and language acqui-	sev, Alexey Skrynnik, Artem Zholus, Aleksandr	606
549	sition: A pragmatic approach. <i>Language Learning</i>	Panov, Kavya Srinet, et al. 2021. Neurips 2021	607
550	<i>and Development</i> , 14(3):170–185.	competition iglu: Interactive grounded language un-	608
		derstanding in a collaborative environment. <i>arXiv</i>	609
551	David A Cohn, Zoubin Ghahramani, and Michael I	<i>preprint arXiv:2110.06536</i> .	610
552	Jordan. 1996. Active learning with statistical models.	Wouter Kool, Herke Van Hoof, and Max Welling. 2019.	611
553	<i>Journal of artificial intelligence research</i> , 4:129–145.	Stochastic beams and where to find them: The	612
		gumbel-top-k trick for sampling sequences without	613
554	Alejandrina Cristia, Emmanuel Dupoux, Nan Bern-		
555	stein Ratner, and Melanie Soderstrom. 2019. Seg-		
556	mentability differences between child-directed and		
557	adult-directed speech: A systematic test with an eco-		
558	logically valid corpus. <i>Open Mind</i> , 3:13–22.		

614	replacement. In <i>International Conference on Machine Learning</i> , pages 3499–3508. PMLR.	Roi Reichart, Katrin Tomanek, Udo Hahn, and Ari Rapoport. 2008. Multi-task active learning for linguistic annotations. In <i>Proceedings of ACL-08: HLT</i> , pages 861–869.	666
615			667
616	Brenden Lake and Marco Baroni. 2018. Generalization without systematicity: On the compositional skills of sequence-to-sequence recurrent networks. In <i>proceedings of the 35th International Conference on Machine Learning (ICML)</i> , pages 4487–4499.	Paul Rodriguez. 2001. Simple recurrent networks learn context-free and context-sensitive languages by counting. <i>Neural computation</i> , 13(9):2093–118.	670
617			671
618			672
619			
620			
621	Angeliki Lazaridou, Adhiguna Kuncoro, Elena Gribovskaya, Devang Agrawal, Adam Liska, Tayfun Terzi, Mai Gimenez, Cyprien de Masson d’Autume, Sebastian Ruder, Dani Yogatama, et al. 2021. Pitfalls of static language modelling. <i>arXiv preprint arXiv:2102.01951</i> .	Paul Rodriguez, Janet Wiles, and Jeffrey L Elman. 1999. A recurrent neural network that learns to count. <i>Connection Science</i> , 11(1):5–40.	673
622			674
623			675
624			
625			
626			
627	Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In <i>Text summarization branches out</i> , pages 74–81.	Diana Rodríguez Luna, Edoardo Maria Ponti, Dieuwke Hupkes, and Elia Bruni. 2020. Internal and external pressures on language emergence: least effort, object constancy and frequency. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 4428–4437, Online. Association for Computational Linguistics.	676
628			677
629			678
630			679
631	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> .	David Saxton, Edward Grefenstette, Felix Hill, and Pushmeet Kohli. 2019. Analysing mathematical reasoning abilities of neural models. In <i>Proceedings of the 7th International Conference on Learning Representations (ICLR)</i> .	680
632			681
633			682
634			683
635	Tambet Matiisen, Avital Oliver, Taco Cohen, and John Schulman. 2019. Teacher–student curriculum learning. <i>IEEE transactions on neural networks and learning systems</i> , 31(9):3732–3740.	Laura A Shneidman and Susan Goldin-Meadow. 2012. Language input and acquisition in a mayan village: How important is directed speech? <i>Developmental science</i> , 15(5):659–673.	684
636			685
637			686
638			687
639	Michael McCloskey and Neal J Cohen. 1989. Catastrophic interference in connectionist networks: The sequential learning problem. In <i>Psychology of learning and motivation</i> , volume 24, pages 109–165. Elsevier.	Katharina Stein, Leonie Harter, and Luisa Geiger. 2021. Shapelurn: An interactive language learning game with logical inference. In <i>Proceedings of the First Workshop on Interactive Learning for Natural Language Processing</i> , pages 16–24.	688
640			689
641			690
642			691
643			692
644	Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. 2016. Asynchronous methods for deep reinforcement learning. In <i>International conference on machine learning</i> , pages 1928–1937. PMLR.	Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-sne. <i>Journal of machine learning research</i> , 9(11).	693
645			694
646			695
647			696
648			
649			
650	Shashi Narayan, Shay B Cohen, and Mirella Lapata. 2018. Ranking sentences for extractive summarization with reinforcement learning. <i>arXiv preprint arXiv:1802.08636</i> .	Oskar van der Wal, Silvan de Boer, Elia Bruni, and Dieuwke Hupkes. 2020. The grammar of emergent languages. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 3339–3359, Online. Association for Computational Linguistics.	700
651			701
652			702
653			703
654			704
655			705
656			
657			
658	Mitja Nikolaus and Abdellah Fourtassi. 2021. Modeling the interaction between perception-based and production-based learning in children’s early acquisition of semantic knowledge.	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In <i>Advances in neural information processing systems</i> , pages 5998–6008.	706
659			707
660			708
661			709
662			710
663	Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In <i>Proceedings of NAACL-HLT 2019: Demonstrations</i> .	Tim Vieira. 2014. Gumbel-max trick and weighted reservoir sampling.	711
664			712
665			
666	J. Rapin and O. Teytaud. 2018. Nevergrad - A gradient-free optimization platform. <a href="https://GitHub.com/FacebookResearch/Nevergrad">https://GitHub.com/FacebookResearch/Nevergrad</a> .	Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohananey, Wei Peng, Sheng-Fu Wang, and Samuel R Bowman. 2020. Blimp: The benchmark of linguistic minimal pairs for english. <i>Transactions of the Association for Computational Linguistics</i> , 8:377–392.	713
			714
			715
			716
			717

Janet Wiles and Jeff Elman. 1995. Learning to count without a counter: A case study of dynamics and activation landscapes in recurrent networks. In *Proceedings of the seventeenth annual conference of the cognitive science society*, s 482, page 487. Erlbaum Hillsdale, NJ.

Ronald J Williams. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256.

Benfeng Xu, Licheng Zhang, Zhendong Mao, Quan Wang, Hongtao Xie, and Yongdong Zhang. 2020. Curriculum learning for natural language understanding. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6095–6104.

Wei Zhang, Wei Wei, Wen Wang, Lingling Jin, and Zheng Cao. 2021. Reducing bert computation by padding removal and curriculum learning. In *2021 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*, pages 90–92. IEEE.

## A Computing the Probability of a Top-K Sample

Our objective is to find the (log) probability of sampling the subset  $(i_1, \dots, i_K)$  from  $\{1, \dots, N\}$  *without* replacement from the categorical probability  $(p_1, \dots, p_N)$ .

Let us first consider sampling  $K$  elements from the  $\{1, \dots, N\}$  *with* replacement. In that case

$$p(i_1, \dots, i_K) = \prod_{k=1}^K p_{i_k}. \quad (1)$$

If we allow for all possible permutations of observing  $(i_1, \dots, i_K)$  we get

$$p(i_1, \dots, i_K) = C \prod_{k=1}^K p_{i_k}, \quad (2)$$

where  $C = K!$ .

To go from sampling *with* replacement, to sampling *without* replacement, we consider event  $A =$  “all sampled elements  $(i_1, \dots, i_K)$  are unique”. Then

$$\begin{aligned} p_{\text{w/o replacement}}(i_1, \dots, i_K) = \\ p_{\text{w/ replacement}}(i_1, \dots, i_K | A). \end{aligned} \quad (3)$$

Applying Bayes Rule gives us:

$$\begin{aligned} p_{\text{w/o replacement}}(i_1, \dots, i_K) = \\ \frac{p_{\text{w/ replacement}}(A | i_1, \dots, i_K) p_{\text{w/ replacement}}(i_1, \dots, i_K)}{p_{\text{w/ replacement}}(A)}. \end{aligned} \quad (4)$$

As in our case all samples in  $(i_1, \dots, i_K)$  are unique we know that

$$p_{\text{w/ replacement}}(A | i_1, \dots, i_K) = 1. \quad (5)$$

Combining this with Equation 2 gives us

$$p_{\text{w/o replacement}}(i_1, \dots, i_K) = \frac{C \prod_{k=1}^K p_{i_k}}{p(A)}, \quad (6)$$

and thus

$$p_{\text{w/o replacement}}(i_1, \dots, i_K) \propto \prod_{k=1}^K p_{i_k}, \quad (7)$$

and

$$\log p_{\text{w/o replacement}}(i_1, \dots, i_K) = \sum_{k=1}^K \log p_{i_k}. \quad (8)$$

From an implementation perspective this this boils down to the following steps:

1. We compute the scores per sentence.
2. We sample  $K$  sentences without replacement, using the GumbelTopK trick.
3. We compute the log probabilities for each score:  $\log \text{softmax}(\text{scores})$ .
4. We compute the log probability of our sample by adding the log probabilities of the elements in our sample, according to Equation 8.

### A.1 Comparison to Prior Work

Our problem of sampling  $K$  sentences as a single action is similar to the problem formulation of using Reinforcement Learning for extractive summarization to optimize for Rouge (Lin, 2004) directly. In this setting  $K$  sentences need to be selected from a document. This results in a very large search space. Narayan et al. (2018) limit the search space by first selecting  $n$  sentences that have a high Rouge score. Then all possible summaries are made with these  $n$  sentences. These summaries are ranked according to their Rouge scores and the top  $K$  sentences are taken as action. This approach

793 has the disadvantage that it limits the search space  
794 heuristically, which does not guarantee that the best  
795 summary is found. [Dong et al. \(2018\)](#) frame the  
796 problem as a contextual bandit problem, which al-  
797 lows them to sample from the true action space.  
798 We choose our approach as it is intuitive, simple  
799 and effective.

## B Additional Results Baseline Experiments Task 1

800

In Table 4 we present the results for our baseline runs on all five seeds.

801

Baseline	Seed	Avg Perplexity	Avg train from test	Avg unigram overlap	Avg bigram overlap	Avg trigram overlap
<i>Random</i>	6639	$193.9 \pm 100.3$	$0.46 \pm 0.14$	$0.46 \pm 0.14$	$0.278 \pm 0.07$	$0.023 \pm 0.009$
	7519	$683.1 \pm 634.3$	$0.52 \pm 0.15$	$0.52 \pm 0.15$	$0.291 \pm 0.10$	$0.030 \pm 0.010$
	1007	$742.5 \pm 159.8$	$0.50 \pm 0.17$	$0.50 \pm 0.17$	$0.298 \pm 0.10$	$0.035 \pm 0.014$
	4520	$160.9 \pm 217.7$	$0.54 \pm 0.16$	$0.54 \pm 0.16$	$0.327 \pm 0.09$	$0.035 \pm 0.025$
	4527	$307.1 \pm 295.1$	$0.58 \pm 0.17$	$0.58 \pm 0.17$	$0.349 \pm 0.10$	$0.035 \pm 0.014$
<i>Oracle</i>	6639	$14.99 \pm 5.364$	$1.00 \pm 0.00$	$1.00 \pm 0.00$	$0.551 \pm 0.06$	$0.072 \pm 0.029$
	7519	$44.37 \pm 58.94$	$1.00 \pm 0.00$	$1.00 \pm 0.00$	$0.611 \pm 0.02$	$0.085 \pm 0.017$
	1007	$68.95 \pm 87.49$	$1.00 \pm 0.00$	$1.00 \pm 0.00$	$0.598 \pm 0.02$	$0.077 \pm 0.025$
	4520	$15.65 \pm 4.616$	$1.00 \pm 0.00$	$1.00 \pm 0.00$	$0.578 \pm 0.02$	$0.087 \pm 0.028$
	4527	$23.66 \pm 21.44$	$1.00 \pm 0.00$	$1.00 \pm 0.00$	$0.624 \pm 0.02$	$0.095 \pm 0.019$
<i>Worst case</i>	6639	$8.46e4 \pm 4.69e4$	$0.00 \pm 0.00$	$0.00 \pm 0.00$	$0.00 \pm 0.00$	$0.00 \pm 0.00$
	7519	$7.03e4 \pm 3.73e4$	$0.00 \pm 0.00$	$0.00 \pm 0.00$	$0.00 \pm 0.00$	$0.00 \pm 0.00$
	1007	$8.17e4 \pm 4.26e4$	$0.00 \pm 0.00$	$0.00 \pm 0.00$	$0.00 \pm 0.00$	$0.00 \pm 0.00$
	4520	$4.78e4 \pm 2.67e4$	$0.00 \pm 0.00$	$0.00 \pm 0.00$	$0.00 \pm 0.00$	$0.00 \pm 0.00$
	4527	$6.69e4 \pm 1.98e4$	$0.00 \pm 0.00$	$0.00 \pm 0.00$	$0.00 \pm 0.00$	$0.00 \pm 0.00$

Table 4: Baseline results for Task 1. Different domains. Averages and standard deviations reported based on five runs per seed.

## C Additional Results Baseline Experiments Task 2

802

In Table 5 we present the results for our baseline runs on all five seeds.

803

## D Additional Results Task 1

804

In this section we present the plots for the  $n$ -gram overlap for Task 1 in Figures 6 and 7.

805

## E Additional Results Task 2

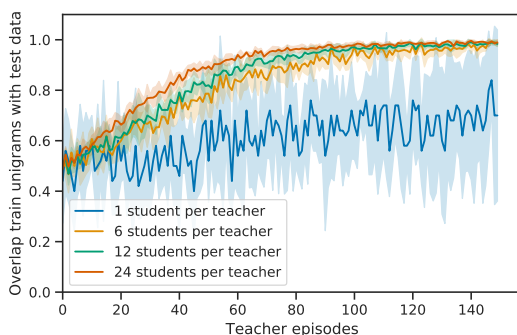
806

In this section we present the plots for the  $n$ -gram overlap for Task 2 in Figure 8.

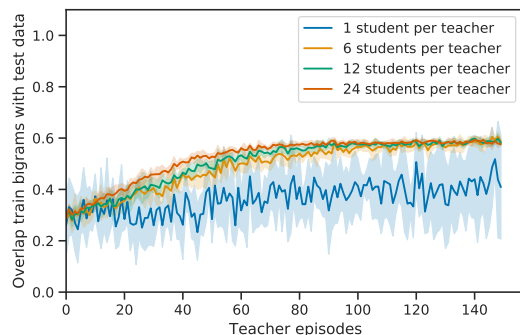
807

Baseline	Seed	Avg Perplexity	Avg train from test	Avg unigram overlap	Avg bigram overlap	Avg trigram overlap
<i>Random</i>	6639	119.0 ± 56.48	0.18 ± 0.04	1.00 ± 0.00	0.401 ± 0.033	0.030 ± 0.020
	7519	162.8 ± 201.9	0.24 ± 0.05	1.00 ± 0.00	0.408 ± 0.044	0.035 ± 0.038
	1007	234.1 ± 192.0	0.24 ± 0.12	1.00 ± 0.00	0.414 ± 0.034	0.034 ± 0.020
	4520	161.7 ± 190.6	0.22 ± 0.04	1.00 ± 0.00	0.410 ± 0.023	0.038 ± 0.033
	4527	342.1 ± 241.4	0.12 ± 0.08	1.00 ± 0.00	0.348 ± 0.024	0.013 ± 0.017
<i>Oracle</i>	6639	6.973 ± 1.534	1.00 ± 0.00	1.00 ± 0.00	0.720 ± 0.044	0.151 ± 0.022
	7519	7.626 ± 2.298	1.00 ± 0.00	1.00 ± 0.00	0.682 ± 0.056	0.177 ± 0.033
	1007	7.895 ± 1.106	1.00 ± 0.00	1.00 ± 0.00	0.726 ± 0.045	0.207 ± 0.025
	4520	6.821 ± 0.619	1.00 ± 0.00	1.00 ± 0.00	0.740 ± 0.073	0.197 ± 0.054
	4527	9.431 ± 3.057	1.00 ± 0.00	1.00 ± 0.00	0.700 ± 0.056	0.174 ± 0.017
<i>Worst case</i>	6639	595.3 ± 297.9	0.00 ± 0.00	1.00 ± 0.00	0.326 ± 0.026	0.00 ± 0.00
	7519	317.2 ± 235.8	0.00 ± 0.00	1.00 ± 0.00	0.311 ± 0.018	0.00 ± 0.00
	1007	508.1 ± 155.7	0.00 ± 0.00	1.00 ± 0.00	0.345 ± 0.017	0.00 ± 0.00
	4520	299.6 ± 124.2	0.00 ± 0.00	1.00 ± 0.00	0.310 ± 0.027	0.00 ± 0.00
	4527	432.8 ± 72.05	0.00 ± 0.00	1.00 ± 0.00	0.330 ± 0.035	0.00 ± 0.00

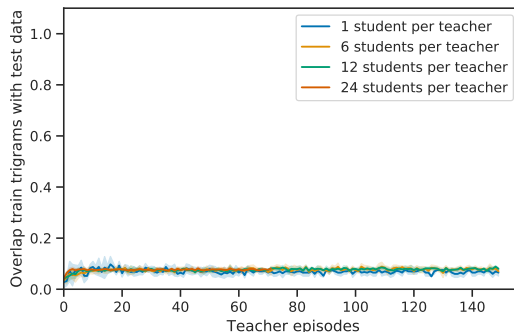
Table 5: Baseline results for Task 2. Different structures. Averages and standard deviations reported based on five runs per seed.



(a) Unigram overlap between train and test data.

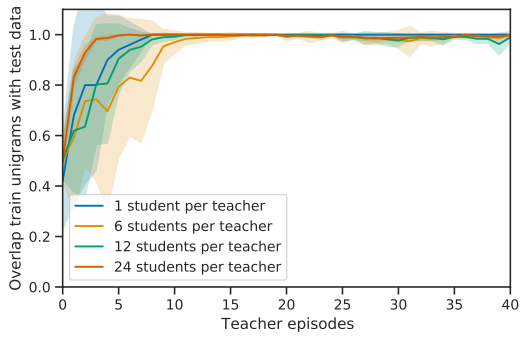


(b) Bigram overlap between train and test data.

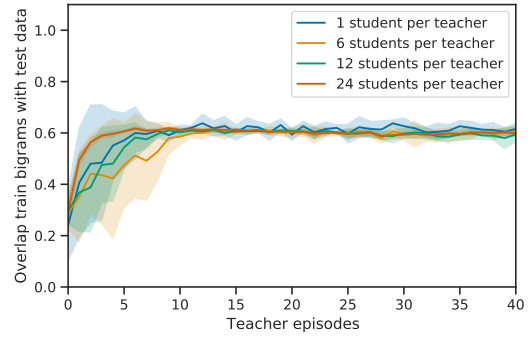


(c) Trigram overlap between train and test data.

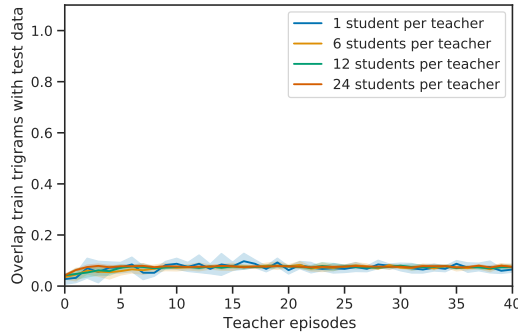
Figure 6: Additional results Task 1 – Different domains. Plots for different numbers of students per teacher. Results per setting reported as average and standard deviation over five random seeds. Average word embedding as sentence embeddings.



(a) Unigram overlap between train and test data.

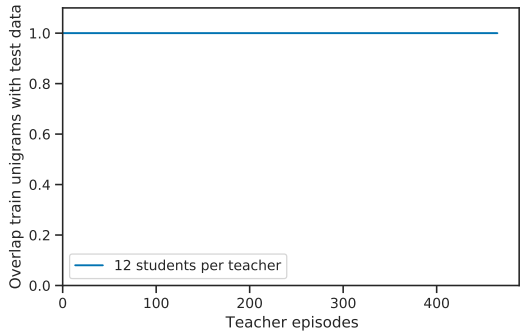


(b) Bigram overlap between train and test data.

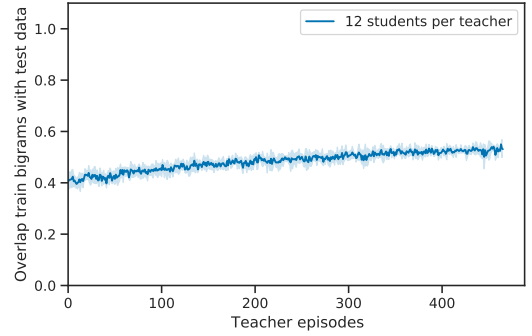


(c) Trigram overlap between train and test data.

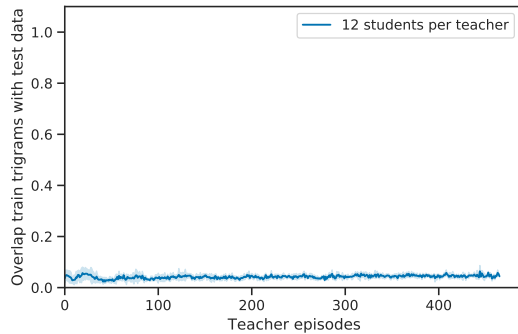
Figure 7: Additional results Task 1 – Different domains. Plots for different numbers of students per teacher. Results per setting reported as average and standard deviation over five random seeds. Average hidden layer embedding as sentence embeddings.



(a) Unigram overlap between train and test data.



(b) Bigram overlap between train and test data.



(c) Trigram overlap between train and test data.

Figure 8: Additional results Task 2 – Different structures. Results per setting reported as average and standard deviation over five random seeds.