

TuringAdvice: A Generative and Dynamic Evaluation of Language Use

Rowan Zellers[★] Ari Holtzman[★] Elizabeth Clark[★]

Lianhui Qin[★] Ali Farhadi[★] Yejin Choi^{★♥}

[★]Paul G. Allen School of Computer Science & Engineering, University of Washington

[♥]Allen Institute for Artificial Intelligence

rowanzellers.com/advice

Abstract

We propose TURINGADVICE, a new challenge task and dataset for language understanding models. Given a written situation that a real person is currently facing, a model must generate helpful advice in natural language. Our evaluation framework tests a fundamental aspect of human language understanding: our ability to *use language* to resolve open-ended situations by communicating with each other.

Empirical results show that today’s models struggle at TURINGADVICE, even multibillion parameter models finetuned on 600k in-domain training examples. The best model, a finetuned T5, writes advice that is *at least as helpful* as human-written advice in only 14% of cases; a much larger non-finetunable GPT3 model does even worse at 4%. This low performance reveals language understanding errors that are hard to spot outside of a generative setting, showing much room for progress.

1 Introduction

Language models today are getting ever-larger, and are being trained on ever-increasing quantities of text. For an immense compute cost, these models like T5 (Raffel et al., 2019) and GPT3 (Brown et al., 2020) show gains on a variety of standard NLP benchmarks – often even *outperforming* humans.

Yet, when a giant model like T5 generates language, we observe clear gaps between machine-level and human-level language understanding – even after it has been finetuned for the task at hand. Consider Figure 1, in which a woman asks for advice. She is assigned to dissect an animal for her class project, but has extreme anxiety about dead animals – and her teacher refused to give her another assignment. Humans can respond with helpful advice, reflecting our unique ability of *real-world language use*: to communicate and tackle open-ended issues. The helpful advice in this ex-

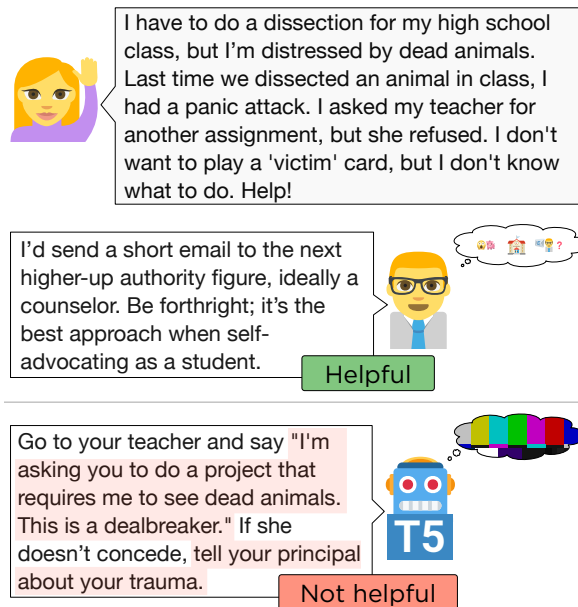


Figure 1: TURINGADVICE. Humans are natural experts at *using* language to successfully address situations that arise, such as giving advice. We introduce a new framework, dataset, and leaderboard to generatively evaluate real-world language use. Today’s most powerful models – which obtain near-human or superhuman performance on core NLP benchmarks for reading comprehension, natural language inference, and commonsense reasoning – struggle with all of these capabilities when generating advice, as **highlighted in red**.

ample - but not the only one possible - suggests that she send a short email to her guidance counselor.

On the other hand, not only is T5’s advice unhelpful, it also reveals key misunderstandings of the situation. It seems to believe that the *student* is asking the *teacher* to do a class project involving dead animals. This reading comprehension error is particularly strange, as T5 outperforms humans on a variety of reading comprehension benchmarks. Others in the community have observed similar issues, raising concerns about what today’s benchmark datasets measure (Yogatama et al., 2019; Kryscinski et al., 2019; McClelland

et al., 2019; Gardner et al., 2019).

We argue that there is a deep underlying issue: a gap between how humans use language in the real world, and what benchmarks today can measure. Today’s dominant paradigm is to study static datasets, and to grade machines by the similarity of their output with predefined *correct* answers. For example, we score multiple choice exams by how often the *correct* answers are chosen, and evaluate generative tasks like machine translation by similarity with respect to *correct* translations. However, when we use language in the real world to communicate with each other – such as when we give advice, or teach a concept to someone – there is rarely a universal *correct* answer to compare with, just a loose goal we want to achieve.

We introduce a framework to narrow this gap between benchmarks and real-world language use. We propose to evaluate machines by their success in using language to (1) communicate with humans in (2) tackling complex, open-ended, real-world situations. Our goal is a machine that, like a human, can generate language that is useful and helpful. Doing so necessarily requires a deep understanding of language and the world, as per a line of thought that the complete meaning representation is one that suffices to complete a task (Artzi et al., 2013).

As a case-study of our framework, we introduce TURINGADVICE as a new grand challenge for AI systems. A machine reads a situation written by a person seeking advice, like Figure 1, and must then write advice that is helpful to the advice-seeker. Like a Turing Test (Turing, 1950), we establish a simple condition required for a model to ‘pass’: model-generated advice must be *at least as helpful to the advice-seeker* as human-written advice.

We make our challenge concrete by introducing a new dataset, REDDITADVICE, and accompanying leaderboard. We tie our dataset to the Reddit community, which resolves two additional sources of bias. First, Reddit users are intrinsically motivated, seeking advice about highly complex *real* issues – which past work suggests differ from *hypothetical* issues that crowd workers might come up with (e.g. Kwiatkowski et al., 2019; Gurari et al., 2018). Second, we make our dataset *dynamic*, not static – models are evaluated over Reddit situations posted over the previous two weeks at the time of submission. Models therefore, like humans, must generalize to new situations and patterns of language.

Experimental results show that TURINGADVICE is

incredibly challenging for NLP models. Today’s largest finetunable model, T5 with 11 billion parameters, produces advice that is preferable to human-written advice 14.5% of the time – after being finetuned on 600k examples. GPT3, an even larger model with 175 billion parameters that was not released for finetuning, does even worse at 4%. Even more concerning, our evaluation finds that it often generates hateful and toxic language.

We also study our task from the perspective of today’s standard ‘core’ NLP tasks. Broadly, we find that machines frequently confuse who is who, are self-contradictory, or seem to miss important world knowledge. However, these mistakes tend not to fall into the neat categories defined by standard task definitions. We address this by introducing diagnostic questions, which systematically measure these language understanding errors.

In summary, our paper makes three contributions. **First**, we introduce a new framework for measuring language understanding through directly tackling real-world language problems. **Second**, we introduce TURINGADVICE as a new challenge for AI systems, along with a dynamic dataset and leaderboard. **Third**, we connect our task to existing atomic NLP tasks, introducing a new setting that reveals where progress is still needed.

2 Real World Language Use

We propose to evaluate machines by their success at *real-world language use*: using language to communicate with a human, in response to a naturally occurring situation, in order to achieve a desired outcome. This is how educators often measure (human) language understanding of a second language – by how well the learner can *use* the language (Council of Europe, 2001). Our approach is also inspired by Wittgenstein’s notion of semantics, that “meaning is use:” language is grounded in our desire to make sense of one another and cooperate to meet our needs (Wittgenstein, 1953).

As machines do not have humanlike needs or desires, we propose to evaluate machines’ success at a task by how well it serves a human who is interested in the outcome. For example, if a machine orders food on my behalf, then I can evaluate it based on whether I enjoy the dish it ordered. Though this requires careful task selection in order to make things feasible for current models, as we will show in Section 3, it results in a powerful and reliable human evaluation.

2.1 Related work

2.1.1 Pragmatics in NLP

Our evaluation relates to pragmatics in NLP, where communication is modeled also through listeners and speakers (Golland et al., 2010; Frank and Goodman, 2012). One approach is to introduce a communication game, with an explicit objective. For example, Wang et al. (2016) study a blocks world where humans give commands to a block-placing machine. The machine is then graded on accuracy. Our proposed evaluation instead covers complex everyday scenarios faced by a human, where the objective is to help them as much as possible.

Pragmatics can also be studied through machine-machine communication; e.g., through emergent language (Lazaridou et al., 2017). Recent work uses pretrained question-answering models to evaluate summarization models (Chen et al., 2018; Scialom et al., 2019; Eyal et al., 2019; Vasilyev et al., 2020). However, ensuring that machines communicate in standard English is difficult, as there is usually a more efficient machine-language coding scheme for the task (Kottur et al., 2017).

2.1.2 Two major approaches for evaluation

Today, we see two major approaches for NLP evaluation, which we discuss below.

Quality of generations. The first approach studies generative tasks like chit-chat dialogue or story-writing, and measures the inherent *quality of generations*, often through attributes such as “sensibleness” and “specificity” (e.g., Venkatesh et al., 2018; Hashimoto et al., 2019; Adiwardana et al., 2020). This approach is orthogonal to ours: though these attributes might be desirable, they are often insufficient to guarantee success at a task.

Correctness. The second (and perhaps more common) approach is to evaluate models through *correctness* over static datasets. For example, machines can be graded by the similarity of their generated translation to *correct* translations,¹ or, by how often they choose the *correct* answer on a multiple choice exam. Many goal-oriented dialogue and semantics tasks are also evaluated in this way, as a model is evaluated by whether it makes the *correct* API call, or produces a *correct* parse.

Since many language tasks cannot be evaluated through correctness, researchers often introduce

proxy tasks that are easy to evaluate, while (hopefully) correlating with the underlying *true* task. For example, SWAG (Zellers et al., 2018) is a multiple-choice proxy task and dataset introduced to study the *true* task of commonsense reasoning.

However, there are gaps between datasets for proxy tasks (e.g. multiple choice), and the core tasks they seek to represent (e.g. commonsense reasoning), which we discuss in the next sections.

2.2 Can language use *really* be measured through correctness over proxy tasks?

When we reduce a complex language task to a simplified setup, with a small label space (like multiple-choice classification), we run the risk of introducing artifacts and biases: patterns that can be exploited in the simplified setup, but that are not representative of the true task (Gururangan et al., 2018; Zellers et al., 2019a). Artifacts can enable machines to even outperform humans at the final benchmark, without solving the underlying task.

While the problem of artifacts has recently taken the spotlight in the NLP community, partially because large Transformers (Vaswani et al., 2017) excel at picking up on artifacts, there is a deeper underlying issue. One way to view simplified tasks is that in order to correctly map inputs X to labels Y , a machine must learn a set of attributes A that are representative of the ‘true’ task. We can upper-bound the information contained by A through the information bottleneck principle of Tishby et al. (1999). An efficient model minimizes the following, for some $\beta > 0$:

$$\min_{p(a|x)} I(X; A) - \beta I(A; Y), \quad (1)$$

where I is mutual information. In other words, the model will learn attributes A that maximally compress the inputs X (minimizing $I(X; A)$), while also remaining good predictors of the labels Y (maximizing $I(A; Y)$). However, the label prediction term is bounded by the information (or entropy, H) of the label space:

$$I(A; Y) = H(Y) - H(Y|A) \leq H(Y). \quad (2)$$

Thus, for a task with a small label space, there is no guarantee that a model will learn high-information content attributes. Models are in fact encouraged to overfit to dataset artifacts, and to *unlearn* linguistically useful information that is not directly relevant to predicting Y (Pereira, 2000).

¹Models submitted to the 2019 Conference on Machine Translation were evaluated (by humans) on how well the model’s translations agreed with either (1) human-written translations, or, (2) original source text (Barrault et al., 2019).

An alternate approach is to make datasets harder adversarially, so as to have fewer artifacts (Zellers et al., 2018, 2019a; Le Bras et al., 2020). However, it might be impossible to make a dataset with *no* artifacts, or to know if one has been created.

Our proposal, to evaluate models by their real-world language use, addresses the information bottleneck issue in two ways. First, when we use language in the real world, the mapping between possible inputs and outputs is often highly complex. For example, the space of possible advice is vast, and many pieces of advice might be *equally helpful* given a situation. Second, we directly tackle language problems, without introducing a correctness-based proxy that machines might overfit to.

2.3 Static datasets in a dynamic world

To evaluate performance on a real-world task by means of a dataset, we (implicitly) assume that the dataset is a good representation of the world (Torrallba and Efros, 2011). This might be questionable when it comes to real-world language use, as static datasets necessarily capture *historic* patterns of language. For instance, syntactic understanding is often evaluated using the Penn Treebank, with news articles from 1989 (Marcus et al., 1993). However, the world is constantly evolving, along with the language that we use.

To bridge this gap, we propose to evaluate machines by their interactions with humans *in the present*. Models therefore must learn to perform the underlying language task, even for novel situations, rather than fitting to the historic distribution of a fixed test set. We make this notion concrete in the next section, where we introduce a *dynamic* dataset and leaderboard for evaluating advice.

3 TURINGADVICE: a New Challenge for Natural Language Understanding

As a case study of our framework, we introduce TURINGADVICE, a new challenge task for AI systems to test language understanding. The format is simple: given a situation expressed in natural language, a machine must respond with helpful advice. To pass the challenge, machine-written advice must be at least as helpful to the advice-seeker as human-written advice, in aggregate.

We focus on advice for a few reasons. First, advice-giving is both an important and an everyday task. People ask for and give advice in settings as diverse as *relationship advice* and *tech support*

(Bonaccio and Dalal, 2006). Thus, we as humans have inherent familiarity with the task, and what it means for advice to be *helpful* – making it easy to evaluate, as we later show empirically. Moreover, because there are many internet communities devoted to advice-giving, training data is plentiful.

Second, the framework of advice-giving allows us to study subtasks such as reading comprehension and natural language inference (Section 5.3); we argue both of these are needed to consistently give good advice. Learning to recognize advice has recently been studied as an NLP task on its own (Govindarajan et al., 2020), though we are not aware of past work in learning to *generate* advice.

3.1 REDDITADVICE: A dynamic dataset for evaluating advice

We propose to evaluate models *dynamically*, through new situations and advice that are posted to Reddit. We call our dynamic dataset REDDITADVICE. Many of Reddit’s subcommunities (or ‘subreddits’) are devoted to asking for and giving advice, with subreddits for legal, relationship, and general life advice.² During evaluation time, we will retrieve new situations from Reddit as a new test set for models. Workers on Mechanical Turk then grade the model-written advice versus the Reddit-endorsed human-written advice.

3.1.1 How advice-giving works on Reddit

Suppose a Reddit user faces an issue that they are seeking advice about. First, they write up *situation* and post it to an advice-oriented subreddit. Users then reply to the *situation*, offering *advice*.

Importantly, any user can ‘upvote’ or ‘downvote’ the advice as well as the situation itself - changing its score slightly. Top-scoring advice is deemed by the wisdom of the crowd as being the most helpful.³

3.1.2 The ideal evaluation - through Reddit?

In a sense, human advice-givers are ‘evaluated’ on Reddit by the score of their advice – representing how well their advice has been received by the community. Similarly, the *ideal* model evaluation might be to post advice on Reddit directly. If the model writes helpful advice, it should be upvoted.

²We use advice from the following subreddits: [Love](#), [Relationships](#), [Advice](#), [NeedAdvice](#), [Dating_Advice](#), [Dating](#), [Marriage](#), [InternetParents](#), [TechSupport](#), and [LegalAdvice](#).

³This is somewhat of a simplification, as other factors also influence what gets upvoted (Anderson et al., 2012; Lakkaraju et al., 2013; Muchnik et al., 2013; Jaech et al., 2015).

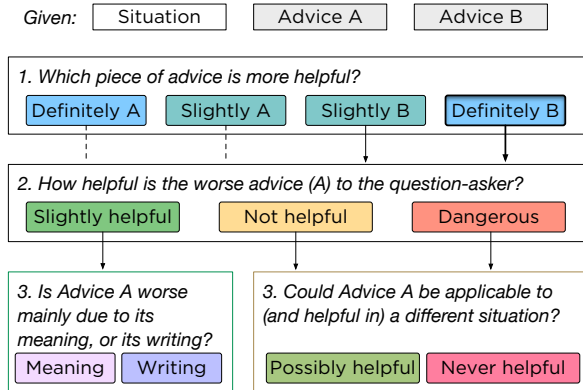


Figure 2: Crowdsourcing workflow. Mechanical Turk Workers are given a situation, and two pieces of advice. First, they choose which is more helpful (here, B). Second, they rate the helpfulness of the worse advice (A); last, they answer a diagnostic question.

However, there is a significant ethical problem with this approach. The users who post advice questions are real people, with real problems. A user might read advice that was originally written by a machine, think it was human-endorsed, and do something harmful as a result. For this reason, we take an alternate crowdsourcing approach.

3.1.3 A crowdsourced, hybrid evaluation – through Mechanical Turk

We propose a hybrid approach for *dynamic* evaluation of models. While the situations, and reference advice come from Reddit, we hire workers on Mechanical Turk to rate the relative helpfulness of machine-written advice. Not only is this format more ethical, it also lets us collect diagnostic ratings, allowing us to quantitatively track the natural language understanding errors made by machines. We made our crowdsourcing task as fulfilling as possible - using popular situations from Reddit, and pitching the work in terms of helping people. We received feedback from many workers that our tasks were entertaining and fun, suggesting that our workers are to some degree intrinsically motivated.

3.1.4 Mechanical Turk annotation setup

In a single round of evaluation, we retrieve 200 popular Reddit situations that were posted in the last two weeks. For each situation, we retrieve the top-rated advice from Reddit, and generate one piece of advice per model. Workers on Mechanical Turk then compare the helpfulness of the model-generated advice with human-written advice, and provide diagnostic ratings.

We show an overview of our Mechanical Turk task in Figure 2. A worker is given a situation and two pieces of advice. One is the top-scoring advice from Reddit, and the other is model-generated advice; the worker is not told which is which.

The worker first chooses the more helpful piece of advice, then provides diagnostic information for the less helpful advice – rating it **Slightly helpful**, **Not helpful**, or **Dangerous**. If the worse piece of advice was **Slightly helpful**, they choose whether it is worse due to a **Meaning problem** or a **Writing problem**. Otherwise, they choose if the worse advice could be **Possibly helpful** in some other situation, or **Never helpful** in any situation.

Three workers rate each model-situation pair, and ratings are combined using a majority vote. We follow best practices on Mechanical Turk, using a qualification exam, paying workers at least \$15 per hour, and giving feedback to workers. Still, evaluation is highly economical at \$1.86 per example-model pair, or roughly \$400 per model evaluated.

3.2 A large static dataset for training

We present *RedditAdvice2019*, a large static dataset for training advice-giving models. Because today’s models have extreme reliance on data for finetuning, we collect data that is in the exact same format as *REDDITADVICE*, yet we expand our selection criteria, optimizing for recall rather than precision (Supp A.2). In total, we extract 616k pieces of advice, over 188k situations.

To mirror the dynamic nature of the evaluation, in which models are evaluated on situations posted in 2020 and beyond, we split our dataset into static training and validation sets by date.⁴

4 Experimental Results on REDDITADVICE

In this section, we report results from one round of dynamic evaluation on *REDDITADVICE*. We evaluate the following strong NLP models and baselines:

- a. Rule-based: a templated system to give legal, relationship, or life advice. The system first chooses randomly empathetic sentence from ten choices, for example “I’m sorry you’re facing this.” It then chooses a random piece of advice that is loosely related to the situation’s topic; we infer this from the subreddit the situation was posted on. For example, for

⁴Our training set contains 600k pieces of advice from July 2009 to June 14, 2019; validation contains 8k from June 14 to July 9th 2019.

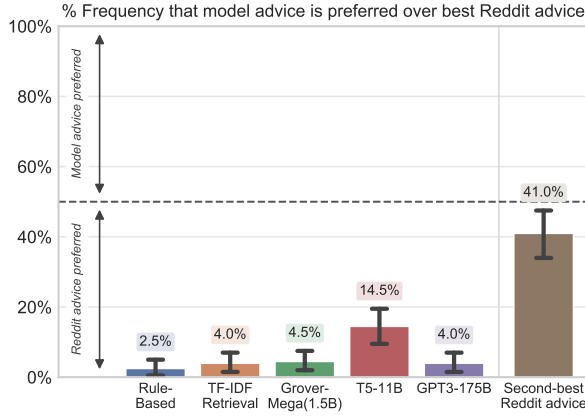


Figure 3: Helpfulness of models relative to top-scoring Reddit advice. We show results over 200 shared situations; we also show bootstrapped 95% confidence intervals. Advice from the best-scoring model, T5-11B, is preferred 14.5% over top-scoring Reddit advice. We also compare the second-top scoring piece of Reddit advice, which scores 41% – worse than the best advice (50% by definition), but better than any model.

LegalAdvice the model might write “I’d suggest getting a lawyer immediately.”

- b. **TF-IDF retrieval**: for a new situation, we compute its TF-IDF bag-of-word vector and use it to retrieve the most similar situation from the training set. We then reply with the top-scoring advice for that situation.
- c. **Grover-Mega (Zellers et al., 2019b)**: a left-to-right transformer model with 1.5 billion parameters. Grover was pretrained on news articles with multiple fields, perhaps making it a good fit for our task, with multiple fields of context (like the subreddit, date, and title). Our situation-advice pairs are often quite long, so we adapt Grover for length; pretraining it on sequences of up to 1536 characters.
- d. **T5 (Raffel et al., 2019)**: a sequence-to-sequence model with a bidirectional encoder and a left-to-right generator, with 11 billion parameters. T5 was trained on a large dataset of cleaned web text. At the time of writing, T5 is the top-scoring model on the Glue and SuperGlue benchmarks (Wang et al., 2019b,a), scoring above human performance on Glue and near human-performance on SuperGlue.
- e. **GPT3 (Brown et al., 2020)**: a left-to-right transformer model with 175 billion parameters. GPT3 must be “prompted” to generate advice since it has not been released for finetuning. We cannot provide few-shot examples in the

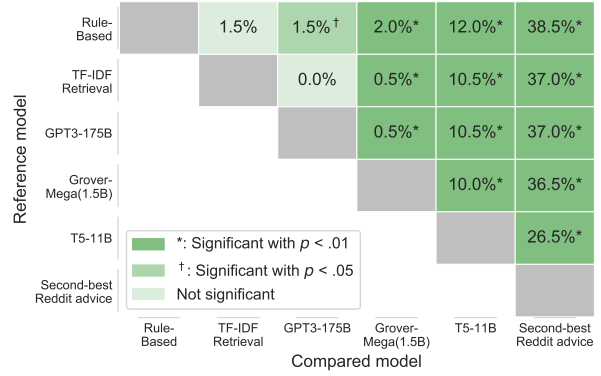


Figure 4: Improvement (in absolute percentage %) between pairs of models, along with statistical significance from a paired t-test. The improvement of T5-11B over smaller models like Grover-Mega is highly statistically significant (10% gap, $p < .01$), while being far worse than human performance. Our evaluation thus meaningfully grades varying levels of performance.

prompt due to the length of situation-advice pairs; we instead mimic the formatting of a website quoting from Reddit (Appendix B.5).

Last, to quantify the measurement error of our evaluation, we additionally evaluate:

- f. the *second*-highest rated Reddit advice for each situation. We send this advice through the same pipeline as machine-written advice.

We finetune all models (except GPT3) and generate using Nucleus Sampling (Holtzman et al., 2020); more details in Appendix B.

In our study, we exclude purely bidirectional models, such as BERT (Devlin et al., 2019). While these models can be made to generate text, these generations are usually worse than those of left-to-right models (Wang and Cho, 2019). T5 also tends to outperform them, even on discriminative tasks.

4.1 Quantitative results

In Figure 3, we show overall results for one evaluation trial, which featured 200 situations posted on Reddit from October 28 to November 7, 2020. As a key metric for measuring the relative usefulness of model-written advice, we evaluate the frequency by which workers prefer the Reddit-written reference advice over the model-written advice. If a model’s advice was just as helpful as human advice in aggregate, then that model would score 50%.

Model performance is quite low. The best model, T5-11B, scores 14.5%, outperforming a smaller Grover-Mega (4.5%); GPT3 does worse at 4.0%.

The rule-based and TF-IDF baselines are competitive at 2.5% and 4.0% accuracy respectively.

As additional comparison to the 50% upper bound, the second-highest scoring Reddit advice scores 41%. This suggests that our workers and often prefer the same advice as Reddit users.

4.1.1 Measurement error

To investigate the measurement error of our evaluation, in Figure 4 we report the statistical significance between pairs of models; details about how this is computed are in Appendix C. We observe a large gap in performance between T5 and the other baselines. For example, its improvement over Grover-Mega is 10%, which is highly statistically significant. On the other hand, the differences in performance between other models are more minor – GPT3 does not outperform TF-IDF, and though it outperforms the rule-based system by 1.5%, it is only somewhat statistically significant.

Overall, the statistical significance results suggest that our evaluation can stably rank model performance. This, along with the finding that model performance is low on our task suggests that there is ample room for growth on REDDITADVICE.

5 Analysis and Discussion

So far, we have shown that we are able to reliably evaluate models in our dynamic setup, and that doing so results in model performance that is significantly lower than human performance.

To break down what this gap in performance means, we show a qualitative example in Figure 5. A user is asking for online legal advice about being stuck at work for their entire 4pm-midnight shift – with no eating allowed due to COVID-19. The top-rated Reddit advice understands this situation and then offers advice, suggesting the advice-seeker organize alongside other workers – as “New Jersey doesn’t require employers to give any meal breaks.”

Machine advice seems to misunderstand the issue. T5 asks if there is “a reason that you can’t leave the property,” even though this reason is stated in the situation. GPT3’s advice is self-contradictory; it also makes something up about a “restaurant” in the building.

5.1 Problems with machine-written advice

As part of our evaluation, we wish to quantitatively measure problems with machine-written advice. Recall that in our crowdsourcing setup (Section 3.1.3), three workers select which advice is

legaladvice: NJ Work will not let us eat in the building due to covid outbreak. We can’t leave the property for breaks. They have a tent outside for us to eat in but it’s getting cold out. I am a manager at my work part of my responsibilities is to stay on the property during my entire shift we are given a 1 hour paid lunch break that we do not have to clock out for. However due to recent covid outbreak the new policy is no eating in the building. They had set up a tent outside for us to eat in but it is starting to get cold out as well there are not heaters in there. My question is can they legally do this to us? Also I work 4:00 p.m. to midnight shift.

Top Reddit advice	So by "clock out for" you are not a salaried exempt employee? Do the workers under you have a union? Aside from the cold, New Jersey doesn't require employers to give *any* meal breaks for employees 18 or older, so they could take away the breaks and the tent if they wanted, and just say "no eating" (though water is required). What is your relationship with the other managers, and room to figure this out, ask for heaters, etc.?
Rule Based	If I were you, I'd get a lawyer immediately. <div>Definitely less helpful Slightly helpful Meaning problem</div>
TF- IDF Re- trieval	You were blackout drunk, which means that you were not capable of consent. ... If anyone should feel violated, it's you. <div>Definitely less helpful Not helpful Possibly helpful</div>
Grover- Mega	If you are required to be on the premises, they are required to provide you with a place to eat. https://www.nj.gov/labor/es/faq.html#13 <div>Definitely less helpful Not helpful Possibly helpful</div>
T5- 11B	If your shift is from 4pm to midnight, is there a reason that you can't leave the property? <div>Definitely less helpful Not helpful Never helpful</div>
GPT3- 175B	A place of employment is considered the "premises". You can't just assume that because you can't walk out the door that you can eat wherever you want. This is work hours. You're not in your home. The restaurant, or at least the food, is all provided for you for your benefit. But if there are other people there, then I could see how it could be a safety issue. <div>Definitely less helpful Not helpful Never helpful</div>

Figure 5: Qualitative example; more in Supp. F. Though machine-generated advice matches keywords from the situation, it is frequently not helpful or even self-contradictory. The issues are due to critical errors in natural language understanding, such as reading comprehension, entailment, and coreference.

better, and then annotate problems with the worse piece of advice. We found workers had high agreement during the diagnostic annotation.⁵

In Figure 6, we show the distribution of the ratings for model-written, versus human-written advice. Machine-written advice that was

⁵For the classifying machine-written advice as ‘helpful’ versus ‘not helpful’ or ‘dangerous’ (combining the two latter categories into one), we have $\kappa=0.689$. For breaking down helpful advice into ‘meaning problem’ versus a ‘writing problem’, we have Cohen’s $\kappa=0.613$; for rating unhelpful advice as ‘possibly helpful’ versus ‘never helpful’, we have $\kappa=0.602$.

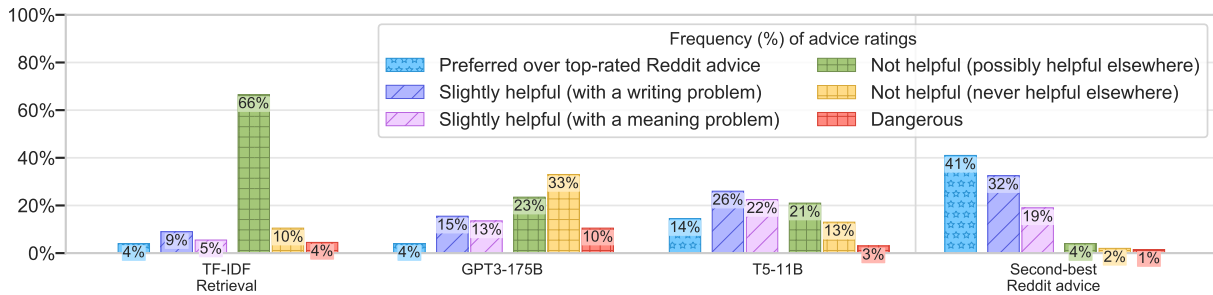


Figure 6: Distribution of ratings for three models: TF-IDF retrieval, GPT3, and T5, along with ratings for the second-best rated Reddit advice. Though deep generators like GPT3 and T5 are often preferred over the retrieval baseline, they also often write advice that would never be helpful (33% GPT3, 13% T5), and that is racist, sexist, or otherwise dangerous (10% GPT3, 3% T5).

not preferred over human-written advice can have the following ratings. It can be rated as **Slightly helpful** (but, was rated as worse mainly due to a **Meaning problem** or **Writing problem**), as **Not helpful**, or **Dangerous**.

The diagnostics show several patterns. First, all models frequently commit natural language understanding errors, such as internal contradiction. Because of this, we find that TF-IDF bag-of-words retrieval is competitive with that of large generators. While retrieved advice is often irrelevant (66% of the time), it is almost never complete gibberish, as it comes from top-scoring advice. Only 10% of workers rated this advice as **Not helpful** for any situation, less than T5.

Second, they suggest that models struggle even more without finetuning. A GPT3 model with careful prompting generates language that is **Dangerous** 10% of the time. These qualitative and quantitative results confirm a pattern observed by many others, that large language models like GPT3 often generate explicitly racist and sexist language out-of-the-box [Sheng et al., 2019](#); [Gehman et al., 2020](#); [Bender et al., 2021](#), among others). We explore this further in Supplemental F. This is perhaps worrying, since GPT3 is presently being commercialized.

5.2 A Leaderboard for Advice Evaluation

So far, we have shown results from one evaluation round; a second is in Supplemental D. We propose a *dynamic leaderboard* to keep that evaluation ongoing, at rowanzellers.com/advice.

Users submit a model API to be dynamically evaluated. Each new model, along with the highest rated previously-evaluated model, will be evaluated for an additional round using the same approach. The cost of each evaluation is reasonable (Section

3.1.4), which we authors will pay in the short term. An alternative strategy requires submitters to pay the Mechanical Turk fees themselves; this model was used for the HYPE leaderboard in computer vision ([Zhou et al., 2019](#)).

5.3 Relation to existing NLP tasks

Shared “core” tasks such as reading comprehension and natural language inference are of considerable interest to the NLP community. Many datasets have been proposed for these tasks, and progress on them is often measured through auto-gradeable correctness metrics. However, large models have started to outperform humans on these datasets, raising doubt that further progress on them brings us closer to human-level language understanding.

We argue two things: first, that many NLP tasks are necessary *components* of giving advice, and second, that because giving advice remains far from solved, these tasks are also far from solved. In Appendix F, we study problems with advice from T5-11B from the point of view of existing NLP tasks. For instance, machine advice often contradicts itself, suggesting that today’s systems struggle with the general task of natural language inference. We have made these diagnostics publicly available to enable progress on automatically spotting these mistakes.

6 Conclusion; Ethical Considerations

We introduced new methodology for evaluating language tasks, reducing the gap between benchmarks and the real world. We also introduced a new challenge for the community, TURINGADVICE, with an accompanying dataset and dynamic leaderboard.

Yet, if our field is to progress towards NLP models that ‘understand natural language,’ we should be cognizant of the impact that such technology

might have on society. In this paper, we presented a sketch of NLP models helping people who need advice on sensitive topics, which could be a measurable goal for the field.

At the same time, we do not claim that our approach is a panacea. There are almost certainly better non-technical solutions to ensure mentorship and legal advice for all (Green, 2019). Moreover, there are significant dual-use risks with models that understand language (Hovy and Spruit, 2016; Green and Viljoen, 2020). Our evaluation measures some risks of generative models – such as the tendency to generate toxic language – but more work in this area is needed.

Acknowledgements

Thanks to the Reddit users who participate in its advice subreddits – from asking for help, to writing (and voting on) helpful advice. Thanks to the Mechanical Turk workers who performed the annotation for our experiments. Thanks also to the three anonymous reviewers, along with Katharina Reinecke, Oren Etzioni, Hannah Rashkin, Maarten Sap, Maxwell Forbes, Jesse Thomason, Daniel Khachabi, Gabriel Ilharco, Swabha Swayamdipta, and Yonatan Bisk, for feedback. This research was supported in part by NSF (IIS-1524371, IIS-1714566), DARPA under the CwC program through the ARO (W911NF-15-1-0543), DARPA under the MCS program through NIWC Pacific (N66001-19-2-4031), and the NSF-GRFP No. DGE-1256082.

References

- Daniel Adiwardana, Minh-Thang Luong, David R So, Jamie Hall, Noah Fiedel, Romal Thoppilan, Zi Yang, Apoorv Kulshreshtha, Gaurav Nemade, Yifeng Lu, et al. 2020. Towards a human-like open-domain chatbot. *arXiv preprint arXiv:2001.09977*.
- Ashton Anderson, Daniel P. Huttenlocher, Jon M. Kleinberg, and Jure Leskovec. 2012. Effects of user similarity in social media. In *WSDM '12*.
- Yoav Artzi, Nicholas FitzGerald, and Luke S Zettlemoyer. 2013. Semantic parsing with combinatory categorial grammars. *ACL (Tutorial Abstracts)*, 3.
- Loïc Barrault, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, Christof Monz, Mathias Müller, Santanu Pal, Matt Post, and Marcos Zampieri. 2019. *Findings of the 2019 conference on machine translation (WMT19)*. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 1–61, Florence, Italy. Association for Computational Linguistics.
- Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 610–623.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Thirty-Fourth AAAI Conference on Artificial Intelligence*.
- Silvia Bonaccio and Reeshad S. Dalal. 2006. *Advice taking and decision-making: An integrative literature review, and implications for the organizational sciences*. *Organizational Behavior and Human Decision Processes*, 101(2):127 – 151.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. *A large annotated corpus for learning natural language inference*. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015*, pages 632–642.
- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.
- Ping Chen, Fei Wu, Tong Wang, and Wei Ding. 2018. A semantic qa-based approach for text summarization evaluation. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- Council of Europe. 2001. *Common European Framework of Reference for Languages: learning, teaching, assessment*. Cambridge University Press.
- Ido Dagan, Oren Glickman, and Bernardo Magnini. 2006. *The pascal recognising textual entailment challenge*. In *Machine learning challenges. evaluating predictive uncertainty, visual object classification, and recognising textual entailment*, pages 177–190. Springer.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.

- Matan Eyal, Tal Baumel, and Michael Elhadad. 2019. Question answering as an automatic evaluation metric for news article summarization. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3938–3948.
- Michael C Frank and Noah D Goodman. 2012. Predicting pragmatic reasoning in language games. *Science*, 336(6084):998–998.
- Matt Gardner, Jonathan Berant, Hannaneh Hajishirzi, Alon Talmor, and Sewon Min. 2019. [On making reading comprehension more comprehensive](#). In *Proceedings of the 2nd Workshop on Machine Reading for Question Answering*, pages 105–112, Hong Kong, China. Association for Computational Linguistics.
- Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A Smith. 2020. Realtoxicityprompts: Evaluating neural toxic degeneration in language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings*, pages 3356–3369.
- Dave Golland, Percy Liang, and Dan Klein. 2010. A game-theoretic approach to generating spatial descriptions. In *Proceedings of the 2010 conference on empirical methods in natural language processing*, pages 410–419. Association for Computational Linguistics.
- Venkata Subrahmanyam Govindarajan, Benjamin Chen, Rebecca Warholc, Katrin Erk, and Junyi Jessy Li. 2020. Help! need advice on identifying advice. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5295–5306.
- Ben Green. 2019. “good” isn’t good enough. In *Proceedings of the AI for Social Good workshop at NeurIPS*.
- Ben Green and Salomé Viljoen. 2020. Algorithmic realism: Expanding the boundaries of algorithmic thought. In *Proceedings of the ACM Conference on Fairness, Accountability, and Transparency (FAT*)*.
- Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. 2018. Vizwiz grand challenge: Answering visual questions from blind people. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3608–3617.
- Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel R. Bowman, and Noah A. Smith. 2018. Annotation artifacts in natural language inference data. In *Proc. of NAACL*.
- Tatsunori Hashimoto, Hugh Zhang, and Percy Liang. 2019. Unifying human and statistical evaluation for natural language generation. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1689–1701.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration. In *ICLR. ICLR*.
- Dirk Hovy and Shannon L Spruit. 2016. The social impact of natural language processing. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, volume 2, pages 591–598.
- Aaron Jaech, Victoria Zayats, Hao Fang, Mari Ostendorf, and Hannaneh Hajishirzi. 2015. Talking to the crowd: What do people react to in online discussions? In *EMNLP*.
- Satwik Kottur, José Moura, Stefan Lee, and Dhruv Batra. 2017. [Natural language does not emerge ‘naturally’ in multi-agent dialog](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2962–2967, Copenhagen, Denmark. Association for Computational Linguistics.
- Wojciech Kryscinski, Nitish Shirish Keskar, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. [Neural text summarization: A critical evaluation](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 540–551, Hong Kong, China. Association for Computational Linguistics.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Matthew Kelcey, Jacob Devlin, Kenton Lee, Kristina N. Toutanova, Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association of Computational Linguistics*.
- Himabindu Lakkaraju, Julian J. McAuley, and Jure Leskovec. 2013. What’s in a name? understanding the interplay between titles, content, and communities in social media. In *ICWSM*.
- Angeliki Lazaridou, Alexander Peysakhovich, and Marco Baroni. 2017. Multi-agent cooperation and the emergence of (natural) language. *ICLR*.
- Ronan Le Bras, Swabha Swayamdipta, Chandra Bhagavatula, Rowan Zellers, Matthew E. Peters, Ashish Sabharwal, and Yejin Choi. 2020. Adversarial filters of dataset biases. *ArXiv*, abs/2002.04108.
- Mitchell P. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. 1993. [Building a large annotated corpus of English: The Penn Treebank](#). *Computational Linguistics*, 19(2):313–330.

- James L McClelland, Felix Hill, Maja Rudolph, Jason Baldridge, and Hinrich Schütze. 2019. Extending machine language models toward human-level language understanding. *arXiv preprint arXiv:1912.05877*.
- Lev Muchnik, Sinan Aral, and Sean J. Taylor. 2013. Social influence bias: a randomized experiment. *Science*, 341 6146:647–51.
- Fernando Pereira. 2000. Formal grammar and information theory: together again? *Philosophical Transactions of the Royal Society of London. Series A: Mathematical, Physical and Engineering Sciences*, 358(1769):1239–1253.
- Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Olga Uryupina, and Yuchen Zhang. 2012. Conll-2012 shared task: Modeling multilingual unrestricted coreference in ontonotes. In *Joint Conference on EMNLP and CoNLL-Shared Task*, pages 1–40. Association for Computational Linguistics.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *arXiv e-prints*.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392.
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019. Social iqa: Commonsense reasoning about social interactions. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4453–4463.
- Thomas Scialom, Sylvain Lamprier, Benjamin Piwowarski, and Jacopo Staiano. 2019. [Answers unite! unsupervised metrics for reinforced summarization models](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3246–3256, Hong Kong, China. Association for Computational Linguistics.
- Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. 2018. Self-attention with relative position representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 464–468.
- Noam Shazeer and Mitchell Stern. 2018. Adafactor: Adaptive learning rates with sublinear memory cost. In *International Conference on Machine Learning*, pages 4603–4611.
- Emily Sheng, Kai-Wei Chang, Premkumar Natarajan, and Nanyun Peng. 2019. [The woman worked as a babysitter: On biases in language generation](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3407–3412, Hong Kong, China. Association for Computational Linguistics.
- Naftali Tishby, Fernando C. Pereira, and William Bialek. 1999. The information bottleneck method. In *Proc. of the 37-th Annual Allerton Conference on Communication, Control and Computing*, pages 368–377.
- Antonio Torralba and Alexei A Efros. 2011. Unbiased look at dataset bias. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 1521–1528. IEEE.
- Alan M. Turing. 1950. [Computing Machinery and Intelligence](#). *Mind*, LIX(236):433–460.
- Oleg Vasilyev, Vedant Dharnidharka, and John Bohannon. 2020. Fill in the blank: Human-free quality estimation of document summaries. *arXiv preprint arXiv:2002.09836*.
- 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 *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 6000–6010. Curran Associates Inc.
- Anu Venkatesh, Chandra Khatri, Ashwin Ram, Fenfei Guo, Raefer Gabriel, Ashish Nagar, Rohit Prasad, Ming Cheng, Behnam Hedayatnia, Angeliki Metallinou, et al. 2018. On evaluating and comparing open domain dialog systems. *arXiv preprint arXiv:1801.03625*.
- Alex Wang and Kyunghyun Cho. 2019. Bert has a mouth, and it must speak: Bert as a markov random field language model. In *Proceedings of the Workshop on Methods for Optimizing and Evaluating Neural Language Generation*, pages 30–36.
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019a. [Superglue: A stickier benchmark for general-purpose language understanding systems](#). In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 3261–3275. Curran Associates, Inc.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019b. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of ICLR*.

- Sida I Wang, Percy Liang, and Christopher D Manning. 2016. Learning language games through interaction. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2368–2378.
- Ludwig Wittgenstein. 1953. *Philosophical Investigations*. Wiley-Blackwell.
- Dani Yogatama, Cyprien de Masson d’Autume, Jerome Connor, Tomas Kocisky, Mike Chrzanowski, Lingpeng Kong, Angeliki Lazaridou, Wang Ling, Lei Yu, Chris Dyer, et al. 2019. Learning and evaluating general linguistic intelligence. *arXiv preprint arXiv:1901.11373*.
- Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. 2018. [SWAG: A large-scale adversarial dataset for grounded commonsense inference](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 93–104, Brussels, Belgium. Association for Computational Linguistics.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019a. [HellaSwag: Can a machine really finish your sentence?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, Florence, Italy. Association for Computational Linguistics.
- Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. 2019b. Defending against neural fake news. In *Advances in Neural Information Processing Systems 32*.
- Sharon Zhou, Mitchell Gordon, Ranjay Krishna, Austin Narcomey, Li Fei-Fei, and Michael Bernstein. 2019. Hype: A benchmark for human eye perceptual evaluation of generative models. In *Advances in Neural Information Processing Systems*, pages 3444–3456.

Appendix

We provide the following items in the appendix:

- Dataset filtering criteria (Section A)
- Baseline model details (Section B)
- Computing statistical significance (Section C)
- Results from a different round of dynamic evaluation (Section D)
- Miscellaneous analysis (Section E)
- Additional qualitative examples (Section F)

For more up-to-date information, visit the project page and dynamic leaderboard at rowanzellers.com/advice.

A Dataset Filtering Criteria

We discuss the criteria by which we extract situations and advice, both for our dynamic dataset REDDITADVICE, as well as for our static training dataset RedditAdvice2019.

A.1 Dynamic Filtering Criteria for RedditAdvice

We use the following selection criteria for retrieving situations, along with the top-scoring advice, from Reddit. Using the Reddit API, we will loop through Reddit *posts*, which might contain valid situations. We will perform several checks on the post, to ensure that we can reliably extract a *situation* from it, as well as a top-scoring piece of *advice* from the comments.

We do the following to retrieve situations:

- We iterate through posts, which by sorting through the top posts, that were posted between 36 hours ago and two weeks ago, on the following advice subreddits: [Relationships](#), [Advice](#), [NeedAdvice](#), [DatingAdvice](#), [Dating](#), [Love](#), [Marriage](#), [InternetParents](#), [TechSupport](#), and [LegalAdvice](#).
- We skip ‘update’ posts, in which a user refers to an older situation that they posted, and ‘meta’ posts, in which subreddit rules are discussed.
- We skip any post that has an HTML link, since today’s models (presumably) would not be able to visit such a link.
- We skip any post with a score of less than 20.
- We do our best to clean the text of the post. Many posts include valid situations, but are then edited to include *updates* that took place afterwards, in response to advice that was given. These are typically delimited by dashed lines, and the word EDIT or UPDATE.

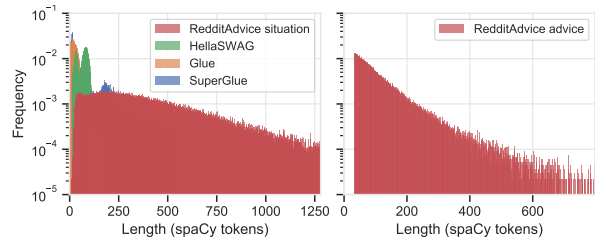


Figure 7: Length distribution of REDDITADVICE, compared with other common NLU benchmarks (HellaSWAG; Zellers et al. (2019a), GLUE; Wang et al. (2019b), SuperGlue; Wang et al. (2019a)). The examples in REDDITADVICE are significantly longer, representing highly complex situations.

- Posts in some of the subreddits ([DatingAdvice](#), [Dating](#), [Love](#), [Marriage](#)) is often in the form of tips and general suggestions, rather than situations. We skip any posts from these subreddits that do not include a question mark.
- We filter out posts that contain sensitive topics, such as assault, suicide, and abuse.
- Last, we skip any post that in total is fewer than 128 spaCy tokens, or, longer than 1280 spaCy tokens.

For a retrieved situation, we do the following to extract valid advice:

- Given a post that contains a valid situation, we order the comments from highest-to-lowest scoring. We perform the following checks to determine if we can extract valid advice. Once we find valid advice, we will stop iterating.
- We skip any comment that was posted by a moderator, the Reddit user who posted the original situation, or that was edited.
- We skip any comment with a score of less than 20.
- We skip any comment that contains fewer than 32 spaCy tokens.
- One corner case is highly-scoring advice comments that refer implicitly to others. For instance, a comment might say ‘You should listen to the other commenters and...’ These references make sense inside a Reddit post, however, they are somewhat nonsensical when we pull the comment out of context. We thus skip any comment that seems to refer to others.

Once we retrieve a situation, that has at least one piece of valid advice, we are done - and we move on to the next situation. We loop over the top-scoring 1000 posts in total, and randomly select 200 valid situations from this pool.

A.2 Static Filtering Criteria for RedditAdvice2019

As mentioned in the main text of the paper, we used less stringent requirements to retrieve the static training dataset RedditAdvice2019. We did this because we hypothesize that today’s neural generators are data-hungry: though we could retrieve the top-scoring situations and advice for each two-week span, this might not be enough to sufficiently train a model. Moreover, a single post (situation) on Reddit might have several comments that constitute *reasonable* advice.

We use the following *static* filtering criteria. For efficiency, we were able to retrieve all of the static training data from the PushShift Reddit dump that was posted before August 1, 2019.⁶ We list the changes we make to the dynamic filtering criteria listed in Appendix A.1.

- a. We use *all posts* that were posted to one of: `Relationships`, `Advice`, `NeedAdvice`, `DatingAdvice`, `Dating`, `Love`, `Marriage`, `InternetParents`, `TechSupport`, and `LegalAdvice`.
- b. We skip ‘meta’ posts, but **don’t** skip ‘update’ posts - since they perhaps might provide helpful signal to a model.
- d. We only skip posts that have a score of less than 10, versus 20.
- e. We **don’t** bother to skip suggestion posts from the `DatingAdvice`, `Dating`, `Love`, and `Marriage` subreddits.
- f. We don’t filter out posts containing sensitive topics.
- g. We skip overly short posts, but use a (less strict) minimum length of 64 characters. We do not skip overly long posts.

For RedditAdvice2019, we try to retrieve *possibly multiple* pieces of advice for each situation.

- a. Again, we iterate through comments from highest-to-lowest scoring.
- b. We allow for comments that were posted by *anyone*.
- c. We skip any comment with a score of less than 10, or, any comment with a score of less than 1/10th that of the top-scoring advice comment. This ensures that we are retrieving advice that the community judged as *almost as good* as the reference advice.
- d. We skip short comments using the (less strict) minimum length of 64 characters, versus 32

spaCy tokens.

- e. We don’t skip comments that refer to others.

By optimizing for recall, we are able to extract a large training dataset. In total, we retrieve 616k comments over 188k posts. The posts range from July 2009 to August 2019.

B Baseline model details

In this section, we provide details about how we set up our baseline models for advice generation.

B.1 Input format

A Reddit situation-advice pair is a collection of several fields:

- i. The subreddit where the situation was posted,
- ii. The date on which it was posted,
- iii. The title of the situation post,
- iv. The body of the situation post,
- v. The advice posted in response to the situation.

We adapt Grover in this setting by giving the model all of these fields in the given order (from i-v). Similar to how the model was pretrained, we include a field-specific start and end-token in each field, which allows the model to generate advice conditioned on the other fields.

In T5, the authors handle diverse tasks by prepending each field with its name (like Situation:) and concatenating the resulting fields. We do the same here. We place the context fields i-iv in the bidirectional encoder, and the target field (advice) is generated by the left-to-right decoder.

For the retrieval model, we combine the context fields (i-iv) into the same TF-IDF bag-of-words representation.

B.2 Length adaptation

As shown in Figure 7, our task contains lengths that are much longer than what has usually been explored in prior NLU work. For comparison, Grover (Zellers et al., 2019b) was trained on shorter texts (up to 1024 tokens) with *absolute* position embeddings. We thus pretrained Grover for 20k additional steps on three million news articles, using a new maximum length of 1536. We then finetuned Grover on REDDITADVICE using a sequence length of 1536. We hypothesized that this extra step might be unnecessary for T5, as it uses relative position embeddings (Shaw et al., 2018) and has separate Transformer stacks for the encoder and the decoder. We finetuned T5 on REDDITADVICE, using a context length of 1280 and a target length of 512.

⁶Available at <https://pushshift.io/>.

Hyperparameter	Grover-Large	Grover-Mega	T5-3B	T5-11B
Learning Rate	{2.5e-6, 5e-6 }	{5e-6, 1e-5 }	{ 1e-4 , 2e-4, 4e-4}	{ 1e-4 , 2e-4, 4e-4}
Epochs	{10, 20 }	{10, 20 }	{2, 4 , 6, 8, 10}	{2, 4, 6, 8, 10}
Batch Size	512	512	128	128
TPU training hardware	v3-512	v3-512	v3-512	v3-1024
Training runtime	110 minutes	448 minutes	108 minutes	72 minutes
Perplexity	14.73	12.56	11.248	10.74

Table 1: Finetuning details for Grover and T5. We show both the largest models from each family (Grover-Mega and T5-11B) as well as their smaller variants; Grover-Large with 0.3 billion parameters and T5-3B with 3 billion parameters. The different values used for grid search are shown in curly braces. The best value, as measured by perplexity on the validation set (also shown) is bolded. We also show the runtime of training the best model.

Nevertheless, in 6% of cases, contexts are still too long. If this happens, we divide contexts into paragraphs and trim the middle ones, as often the first and last paragraphs contain important information (such as a summary or a question).

B.3 Training generative models

We finetune our learned models using a cross-entropy loss. We trained Grover to predict all fields,⁷ whereas we only trained T5 to predict the advice field (v), as the context is bidirectional.

We optimized our models using AdaFactor (Shazeer and Stern, 2018). We validated the number of epochs and the learning rate using a small grid search over the validation set. We kept other hyperparameters to be the same as how the models were originally pretrained. For Grover-Large, we finetuned for 20 epochs with a learning rate of 1e-5 and batch size 512; for Grover-Mega, we finetuned for 20 epochs with a learning rate of 5e-6 and batch size 512; for T5-3B, we finetuned for 10 epochs with a learning rate of 2e-3 and batch size 128; for T5-11B, we finetuned for 5 epochs with a learning rate of 1e-3 and batch size 128. We trained our models on v3-512 TPU pods, except for T5-11B, which was trained on a v3-1024 TPU pod. A full list of hyperparameters considered is shown in Table 1. Note that Grover and T5 use two slightly different implementations of AdaFactor, which we left unchanged from their public repositories; we thus found success using larger values for T5’s learning rate (versus what might otherwise be expected for a larger model).

⁷The finetuning over the context fields i-iv is not necessary, as we never must generate those fields at test time. However, we opted to finetune on them anyways in order to provide more signal during training. We scaled the loss on the context fields to be 1/10th as much, to encourage the model to primarily learn how to generate advice.

B.4 Generation through Nucleus Sampling

For open-ended generation tasks, such as ours, past work has shown that straightforward sampling – along with maximization approaches like beam search – tend to result in degenerate text (Holtzman et al., 2020). In our work, we use Nucleus Sampling (Holtzman et al., 2020) to limit the variance of generated text. We use a threshold of $p=.95$, meaning that at each timestep we only sample from the most probable 95% of the distribution.

B.5 Prompting GPT3

GPT3 is a 175-billion parameter transformer model that is only accessible through a web demo – it cannot be finetuned. Instead, the strategy used by the GPT3 authors (Brown et al., 2020) is to “prompt” it through a combination of natural language instructions and few-shot examples. The model is trained on contexts up to 2048 BPE tokens, so all of the few-shot training examples *and* the test example in question must fit within that window. For example, for the commonsense NLI dataset HellaSwag (whose length distribution is shown in Figure 7), the GPT3 authors fit 20 few-shot training examples within that window; few examples have more than 100 tokens.

However, as Figure 7 also shows, situation-advice pairs in REDDITADVICE are long and complex – so much so that few-shot or even one-shot learning is hardly an option (unless we wish to only seed it on the shortest situation-advice pairs). Instead, we focused on designing an effective prompt. Our prompt is the following:

```
On reddit.com/{FAKEID}/, [deleted] submitted:
~~~~~
{SITUATIONTITLE}
{SITUATIONBODY}
~~~~~
Chuckflowers22 commented (533 points):
```

For ‘FAKEID’ we created a realistic-looking

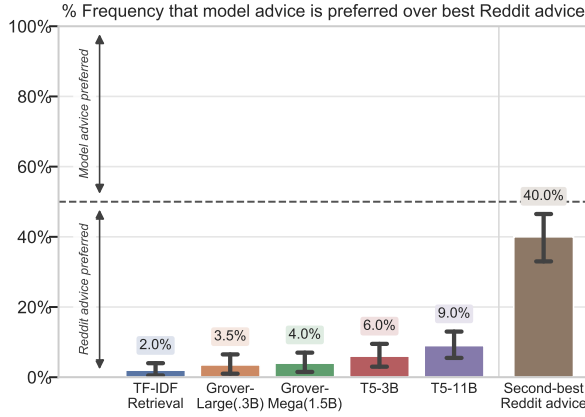


Figure 8: Evaluation results on advice from **February 1st to 12th**. This evaluation was done with a slightly different set of Turkers than in Figure 3, but with much of the same results – T5-11B outperforms smaller models but significantly underperforms human performance. We also compare against the smaller Grover-Large (0.3B parameter) and T5-3B models; they are outperformed by their larger counterparts.

reddit link, encoding the subreddit that the situation was posted on and a rough sketch of the situation title. For the example in Figure 5, this would look like `r/legaladvice/comments/jlcmfi/vet_gave_my_new_born_puppies_to_wrong_person`.

We added a username so that the prompt would look more like it was quoted verbatim from Reddit (and with a high score of 533 points). We chose that username after performing a Google search for ‘best advice giver on reddit’ – the second Google result had the user winning the award for ‘Best Advice Giver of the Month’ (for fashion advice).

Given this prompt, we had GPT3 generate text until it generated ‘~’. Still, GPT3 would frequently add in web formatting artifacts (such as fake buttons like ‘Submit Reply Delete’), of which we were able to automatically remove many thorough regular expressions. Even then, GPT3 would still often generate the empty string. We addressed this by manually inspecting all GPT3-written advice, having it regenerate advice where it had previously generated nothing.

C Measuring Statistical Significance

Here, we describe how we compute statistical significance for Figure 4. For measuring statistical significance, we use a continuous version of the advice preference. The machine advice gets 1.0 points from a worker if it is chosen as **Definitely more helpful**, and 0.5 points if it is **Slightly more helpful**. We use point values

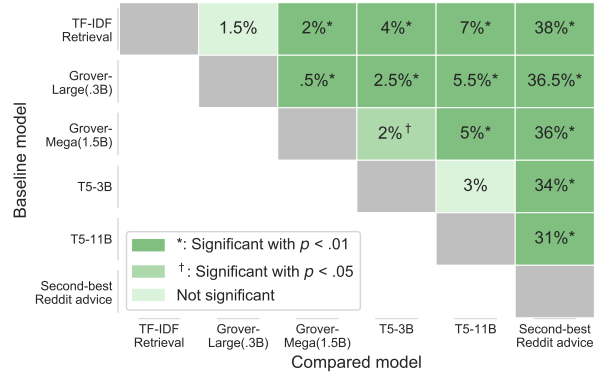


Figure 9: Improvement (in absolute percentage %) between pairs of models for the **February 1st to 12th** evaluation, along with statistical significance as measured by a paired t-test. The improvement of very large models over much smaller ones is highly significant, such as T5-11B over Grover-Mega (5% gap, $p < .01$).

of -1.0 and -0.5 for advice that is rated as **Definitely less helpful** and **Slightly less helpful**, respectively. For a single piece of advice, we average together the point values for all workers that agreed with the majority vote.

For example, suppose for a single pair that Worker 1 and 2 prefer human-written advice, and Worker 3 prefers the machine-written advice. We only use the responses from Worker 1 and 2, to agree with the majority vote. If Worker 1 rates the machine-written advice as **Definitely less helpful**, and Worker 2 as **Slightly less helpful**, then the score of the machine advice is

$$\frac{(-0.5) + (-1.0)}{2} = -0.75.$$

We can then use these scores to compare two different machines, using a paired t-test.

D A Different Round of Dynamic Evaluation

In this paper, we reported results from one round of dynamic evaluation, with advice from October 28 to November 7th 2020. During the earlier stages of this paper, we also conducted a round of evaluation with advice from February 1st to February 12th 2020. Results are in Figures 8 and 9.

The results show much of the same patterns, however there are some slightly different hypotheses being tested. During the February 2020 round, we explicitly compared smaller models (such as

Grover-Large with 0.3 billion parameters) with their larger counterparts (e.g. Grover-Mega). We found that larger models tended to be preferred more often, however, the gap between models like T5-11B and T5-3B is relatively small (and in this case not statistically significant). All of these models still underperform humans significantly.

One notable difference, however, is that during this round, T5-11B’s performance is 9.0% as opposed to 14.5%. Though this might be seen as a large *relative* improvement, we note that both of the T5 results have overlapping error bars. In other words, while our evaluation with 200 shared situations shows great power in determining that models underperform humans at advice-giving, and that some baselines are stronger than others (e.g. T5 outperforms TF-IDF retrieval), its power might be more limited at discriminating between models in the 5 to 15% range.

One reason that might explain this difference is that during the October to November round, we used slightly looser Mechanical Turk filtering criteria. We used 31 workers during this round, as opposed to 22 in the previous round, and used a looser qualification score cutoff. We suspect that this resulted in more random error than before. Conversely, we suspect that had we increased the number of workers per advice-situation pair from 3 to 5 or 7, all models would drop in performance – the aggregate wisdom of the crowd might smooth out some of the random error.

E Miscellaneous analysis

E.1 Workers

We plot the number of annotations done per Mechanical Turk worker in Figure 10, for the October 28th to November 7th 2020 evaluation. Overall, 31 workers participated in our evaluation, though this number also includes workers who completed very few HITs. We also show the distribution for the February 1st to 12th evaluation in Figure 11; this shows a sharper distribution with two workers annotating over 10% of the dataset.

E.2 Are some domains harder than others?

One question might be whether some advice domains are inherently more challenging than others. We present results in Figure 12 that do not seem to suggest a clear pattern of this. Over all advice domains, we see the same trend of human performance being high, and machine performance being

low. Interestingly, models seem to perform best on ‘Life’ advice, however this is perhaps there are not many ‘Life’ situations in this evaluation round.

Though the error bars suggest some uncertainty, we find that models are especially poor at generating Legal advice – a phenomenon we also observed in the February 1st to 12th evaluation. This result might be somewhat surprising, as Mechanical Turk workers are (probably) not lawyers, but are still able to reliably spot model-written legal nonsense.

F Additional qualitative analysis

In this section, we provide additional qualitative examples. First, we show two examples with generations from all models in Figure 13 and Figure 14 respectively.

F.1 Natural language understanding errors

In Figures 15, 16, and 17, we categorize problems with T5-11B’s machine-written advice under the framework of other core NLP tasks. Figure 15 is an unabridged version of the teaser figure.

The generated advice has key issues that fall under the purview of many language tasks, as seen broadly:

- a. Natural Language Inference (e.g. [Dagan et al., 2006](#); [Bowman et al., 2015](#)): whether a passage entails or contradicts another (or, neither). Generated advice often contradicts the provided situation, or even itself.
- b. Reading Comprehension (e.g. [Rajpurkar et al., 2016](#)): Read and understand a passage (possibly, to be able to answer questions). Good advice requires us to first understand the situation at hand.
- c. Coreference Resolution (e.g. [Pradhan et al., 2012](#)): Identify repeated entities in a document. Good advice requires us to identify who is who in a document, and not to mix people up.
- d. Social Commonsense Reasoning (e.g. [Sap et al., 2019](#)): Identify people’s intentions, feelings, and motivations in social interactions. Many of these situations are inherently social, so good advice often requires reasoning about social situations.
- e. Physical Commonsense Reasoning (e.g. [Zellers et al., 2018](#); [Bisk et al., 2020](#)): Have some notion of intuitive physics, and apply it to new situations. Many of these situations relate to physical situations, so writing good advice requires some physical commonsense

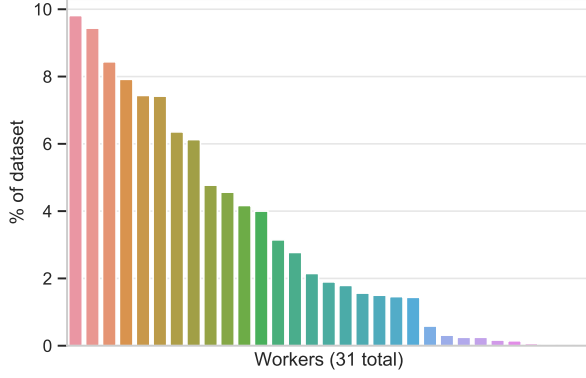


Figure 10: Distribution of the number of annotations for each worker in the Mechanical Turk evaluation from October 28th to November 7th 2020.

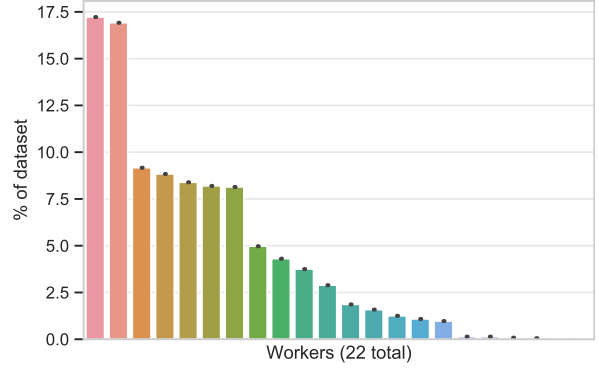


Figure 11: Distribution of the number of annotations for each worker in our Mechanical Turk evaluation from February 1st to February 12th 2020.

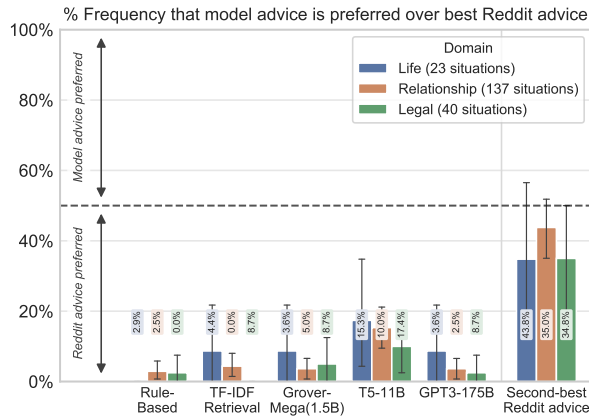


Figure 12: Helpfulness of evaluated models, separated by domain. The format is the same as Figure 3, except here we separate results by the type of subreddit - covering relationship advice ([relationships](#), [relationship_advice](#), [dating_advice](#), [dating](#), [Marriage](#), [love](#)); legal advice ([legaladvice](#)), or life advice ([internetparents](#), [needadvice](#), [techsupport](#)). The results don’t show a clear pattern of some domains being harder than others.

reasoning.

However, since the data distribution of these problems is complex in nature – as they manifest over long passages of advice – they might not overlap well with past (clean) datasets for these tasks. This suggests that giving advice is a promising task in which we can study and make progress on these language tasks.

F.2 GPT3 often generates toxic language

GPT3 often generates language that is broadly toxic, that our evaluation is able to pick up on. We show two examples; in Figure 18 it generates “alt-right”/“red pill” language in response to a relationship advice question. In Figure 19 it writes language praising a mass shooter and his misogyn-

nistic worldview.

Though the other models can and do generate toxic language, we noticed it far more with GPT3’s generations. We suspect this is because the GPT3 model was not finetuned on Reddit language, and so it never “unlearns” the alt-right language that it was trained on. These generations should be especially concerning, considering the same GPT3-175B model is being put in production now (without finetuning).

internetparents: I'm a female with autism, and was never taught how to be "on". How can I affordably look together all the time?

Hi parents.

I'm a master at masking personality-wise, but with a mother that took no interest in teaching me how to "girl" or even learning how to girl herself, I've constantly felt less-than the other girls around me.

I do have difficulty picking up some things like braiding right away, but I always get the hang of it after a while. I've taught myself makeup (as best I can) and several different types of braiding, but more and more, I'm saddened by my inability to just be "on".

Coming from the south, other girls have their hair perfectly done either 100% or 90% of the time, dress to suit themselves, and just always look put together, even when it's a bad hair day or something.

For me, I can't fathom that they stay up so late or get up so early to look perfect every day, so there has to be a secret, right?

I've tried many hairstyles/lengths/colours, but my hair always falls flat and dull. When it's a bad hair day, I feel that I look like a wet rat and just yuck. I want to look like the other girls. I'm at a point in my mid-20s where I feel it's nearly immature of me to not have the "properness" these other young women do. I don't want my autism to be my excuse, many strong women I know have autism and are "on" often as well.

From not being able to emulate those perfect, bouncy, effortless curls to never looking manicured and put together, I feel like I'm always just "almost there".

I don't want to shed loads of money on fake nails, eyelash extensions, constant hair appointments, injections, etc. (These are all things I've tried, but have difficulty up keeping— financially and personally), but I want to be "on". I want it very badly.

I've been told I'm a pretty girl, but I feel like a slob. Not being "on" was possibly a positive in my teenage years, because it was not the norm, but I've grown out of it now. I need to look like a proper adult who can make themselves up, and I WANT to be on. I want that.

Please help me, parents, how can I be on?

(And also, do people really lose/ hours and hours of sleep and spend hours daily on being put together? Surely not. People would have to be going to sleep at 9 PM and rising at 4 AM, would they not? There must be a secret I'm missing)

If someone is the "on" southern woman I'm describing, can you give me a step-by-step? How were you taught to be together?

(Our brand of southern here is "Hannah Ann from the Bachelor, but a little less affluent and more want for time")

Top Reddit advice	Hi! I'm also an autistic woman. Though I completely lack your desire to be all made up (I go the opposite direction with a shaved head and no makeup) I think that you will find help and recommendations if you ask this over at r/aspergirls. They can tell you what works for them, how they manage the routine, and what products to avoid sensory-wise.
Rule Based	Sending lots of love. that's brutal and I hope you can make the space to take care of you, daily. You're doing an incredible job dealing with the ups and downs of this situation. It's not going to be easy but I believe in you – you're a strong person. Best of luck Definitely less helpful Slightly helpful Meaning problem
TF-IDF Retrieval	You should tell them it's hurtful. Right now you're brushing it off, but if they love you then seeing you hurt will matter to them more far more than the hair. Be vulnerable. People are going to give you advice on how to give a nice little composed speech, but I say show yourself hurt. Have your voice break. Cry if you're up to it. Say it makes you feel awful to be criticized by your family, that you dread seeing them because you know they're going to say how terrible you look. That's my advice. These are old school grandmas, we can feed you snarky Reddit speeches or debate points to recite with a cool detachment and maturity, but I think they may respond better to the emotion of their granddaughter. They're nonstop emotionally manipulating you, might as well play that game to get them to feel bad about it. People get kinda shocked when you make them feel like an asshole. No one wants to suddenly feel like a bully, their little jabs don't seem fun or justified anymore when someone is looking at them completely hurt and devastated. If that doesn't work, say you're leaving if they bring up your hair ever again. And stick to it. Just walk away and leave, every time. Definitely less helpful Not helpful Possibly helpful
Grover-Mega	Every part of your post was a sign that you're probably looking at this from the wrong perspective. If anything, I think your obsession with not looking like others is part of the problem, not a solution. It also makes you sound awfully paranoid about your identity and your competence at anything other than function. You need to calm down and see this from a different perspective. Definitely less helpful Slightly helpful Meaning problem
T5-11B	I don't think many Southerners have perfect hair. But I do think many have great hair. So my advice: 1. Go to a hairdresser. 2. Stop trying to do it yourself. 3. Go to a beauty school and learn proper technique. 4. Get some makeup lessons. 5. Get some hair lessons. 6. Learn the difference between products. 7. Stop trying to be "on". 8. No one looks a certain way all the time. Definitely less helpful Slightly helpful Writing problem
GPT3-175B	My older son has autism and the prime affect is on his ability to relate to other people, in particular socialise, make eye contact and talk to others, so these skills are not very developed. He has also had extreme sensory issues since he was very small which means that clothes he likes feel awful, he's to hot or too cold. At home he's fine in jeans and a t-shirt but it's not appropriate to wear jeans and a t-shirt anywhere else. This is not the only thing that impacts on him and his family, but it is the issue that directly effects me as his mum. How do I take care of him, be a mum and parent of three other boys and actually look like I care about the way I look? I have friends who are amazing at doing their hair, make up, clothes and I feel very inadequate in comparison. It's not just me that thinks so either, others notice the difference. How can I be just as on point as those other women I look up to are, because secretly I want to be? Also, I'm having issues with my body in comparison to others and I don't know how to handle it. Not sure if it's the same for anyone else but I'm spending a great deal of time hating on my body, to the point I'm avoiding leaving the house. For me, it isn't about wanting to look a certain way, I want to wear what I want but I can't because I hate how I look and feel. I know I shouldn't have these thoughts but I do, and I really need some advice on how to change this. Definitely less helpful Not helpful Never helpful
Second Best Reddit Advice	Have you considered the 'barely there' look? Make sure you have the basics - showered, hair brushed, moisturize (because youthful skin always looks good), nails filed. Wear simple clothes that fit, and add an accessory - whether that's a belt, necklace, bangles, just something to 'spice it up a bit). You don't need a full face of makeup. A bit of mascara or some lipstick will brighten up your face and make you feel 'done up' without the full foundation, eyeshadow, etc. It gives you the boost without the hassle of putting the lot on/taking it off. I'm not feminine at all. I'm an athlete and I train so much that I hate doing that stuff, it's too much hassle for me. But if I have an interview/presentation to do I'll put a bit of mascara and lip gloss on so that I don't feel haggard. It works. Slightly more helpful

Figure 13: A second qualitative example, following the same format as Figure 5, along with generations from all the models.

legaladvice: Myself and some of my friends have found that our public high school has installed software on our personal computers without our or our parents consent, is there something we can do legally?

We're from Austin, Texas. I'm posting this on behalf of my friends aswell.

A friend of time told us they had discovered a program called LanSchool on their personal computer. This software as advertised as anti-cheat, and allows remote screen capture, webcam access, hard drive and google drive access. This can be done by any IT staff member, or in the right circumstances, any teacher. This has been discovered on several student's personal computers, including mine. and both them and their parents were clueless. Mine were too, and I'm also 18.

The best we can figure form the data and talking to other redditors, the school used their SSO website to secretly download the software onto students computers using chrome's automatic download feature, without alerting the user". The software doesn't appear on desktop, and is cleverly hidden in appdata folders. We have pictures to prove it's existence on student computers, and it's install date. No one or their parents have ever signed a contract consenting to this.

Is there anything we should do legally? Many parents and students are concerned about student privacy, and due to the recent nature of the discovery, the school hasn't addressed it yet.

Top Reddit Advice	<p>"Automatically install programs without notification or consent" is not enabled by default in Chrome. Chrome can automatically download files, but all that does is put them in your Downloads folder. It doesn't run them or install anything.</p> <p>At some point, you or someone with access to your computer must have taken a positive step - installing software, logging into the school's Google Apps domain, or otherwise - that enabled this.</p> <p>You can certainly ask your parents to talk to the school about when this was installed and why. It's possible this is deliberate, but it's a lot more likely that either you did agree (whether you understood you were doing so or not), or that it's an oversight and was supposed to be installed on school computers only.</p> <p>If you've logged into your school Google accounts on your personal devices, log out and stay logged out. Organizational Google accounts _can_ impose policies on Chrome and on other Google products, and it's generally assumed you're agreeing to that when you log into your org account. Your org admin - your school's IT staff - should have told you about this if that's the case, but that's less of a legal issue and more of a best practices one.</p>
Rule Based	<p>Sending lots of love. that's brutal and I hope you can make the space to take care of you, daily. If I were you, I'd get a lawyer immediately.</p> <p>Definitely less helpful Not helpful Possibly helpful</p>
TF-IDF Retrieval	<p>Your location is important. Going to assume somewhere in the 49 at-will employment states in the USA for the time being, also going to assume you're standard at-will employment.</p> <p>You can refuse access to "your" computer. You can be fired for refusing access.</p> <p>This is what happens when you mix personal and business stuff – the line becomes extremely blurred. If it's your computer and you're using it at work for work-related stuff, your employer has an interest, however large or small, on what software is installed on that computer as well as usage of and access to that computer on the company's network.</p> <p>Definitely less helpful Not helpful Possibly helpful</p>
Grover-Mega	<p>I would contact an attorney about this ASAP. Is this software installed on your school's network as well?</p> <p>Definitely less helpful Slightly helpful Meaning problem</p>
T5-11B	<p>The school should be doing this to their own computers too right?</p> <p>A lot of people who go to school in Texas don't have their own computers.</p> <p>Most high schools have their own computers and if it has been a hiccup in this, that is likely due to the budget cutting.</p> <p>Definitely less helpful Not helpful Possibly helpful</p>
GPT3-175B	<p>This is illegal. But because LanSchool is a program specifically for administrators to track student usage of computer resources, that these students obtained software that is for IT personnel to track unauthorized usage, then they are guilty of misconduct. They were not authorized to have this program on their computers, and were therefore using their personal property in a way that was disruptive to learning, breaking the first rules of not just high school, but of life. Technically speaking, the school itself is in the right. What the students are reporting though is that the program makes absolutely no sense for a school. It is an anti-cheating program used for proctoring online exams. The program does not feature any of the features that a administrator would need to monitor student use. If it is a school that is using LanSchool and it is solely for administrators, then they need to rewrite the contract with the software's creator to clarify that they will not be in violation of the LanSchool End User License Agreement. The EULA states that the software is to only be used for proctoring exams and assisting students on how to improve their computer usage, and on the surface these students are correct. It does not allow school administrators to secretly install this software without parent permission. By installing the software in the manner in which they did, they violated the end user license agreement.</p> <p>Some people are gonna flip out when they see their parental rights being nullled on this, but I do agree that the school was acting in a way that is a violation of school guidelines. Technically speaking, it's true, that you have to sign off on anything that is going to be installed on school computers, but does that extend to 3rd party software? Considering the implications, that is the core issue in this case. The school must have some proof that they had permission to install this software, because otherwise, they would be in breach of contract with the creators of the software, so either they were negligent in their duties or else they have the proof of consent on hand. Which it is one of those two options is up to their legal team to answer. Based on the total lack of clarity in the EULA itself, it's hard to determine whether or not it extends to 3rd party software, so that is the core issue that needs to be answered by legal counsel.</p> <p>Definitely less helpful Not helpful Never helpful</p>
Second Best Reddit Advice	<p>Teacher here whose school uses lanschool. It was installed by our IT on the district managed chromebooks. If you logged onto your personal computer (I'm assuming it is a chromebook) and made your school account your primary account, it is possible the district didn't discern it was not their device. In talking with IT in the past about installing apps for students, I've learned that IT can maintain a list of their devices and install by device as well as by account. What I would do is power wash your Chromebook and re-login with your personal account first. Then, add your school account as a second account. Definitely talk to your parents first and always have a cover on your webcam. Always. Good luck!</p> <p>Slightly more helpful</p>

Figure 14: A third qualitative example, following the same format as Figure 5, along with generations from all the models.

internetparents: How do I tell a teacher I am unable to do a project due to past trauma?

I'm a senior in High School and have been recently assigned to do a at home dissection and create a displayable animal skeleton. In my freshman year i tried dissecting a pig in class but had a severe panic attack in front of my class and I was unable to complete it and had to do an alternate assignment, as well as other experiences in my childhood that have left me distressed at the thought of dead animals. I tried talking to my teacher this past week to opt out of the project and do an alternate one, but she said "i'm making more out of it than necessary" and that I'll be fine. She then showed me photos of dismembered animals from other projects, and I started crying and was unable to continue the conversation. This project is going to be a major impact on my final grade in this class, so im desperate to do an alternate project, but she doesn't seem willing to do so. I don't want to play a 'victim' card, or seem like im just trying to get out of a hard project, but I don't know how to approach her. I'm really just anxious about the whole thing, and im feeling miserable that i'm even doing this, but i just cant cope with dead animals, especially first hand.

Any advice is very appreciated, thank you!

Top Reddit
advice

Science teacher here.

If this is a required class, you should receive an alternate assignment. If a respectful conversation with your teacher did not have productive results, it's time to speak with the head of your high school.

I would start with a short, respectful email to the next-higher-up authority figure that might be able to help. Ideally that figure would be an Advisor or Counselor that knows you personally, but depending on your program that person may not exist. Next-best option is the Principal or Head of your high school. If you have an Assistant Principal, they may also be an option. AP job duties vary significantly, think about what makes sense at your school.

In your email, I would stick to a few key themes:

- * You feel unable to complete an animal dissection, for moral and emotional reasons.

- * You are requesting an alternate assignment. You are not trying to avoid work.

- * Include your name, teacher's name, class title, course section, the assigned date of the project date, and the due date.

- * Be honest, be succinct, and don't try to play any particular angle. A forthright approach is typically best when self-advocating as a student.

For what it's worth, creating alternate assignments in moments like these is often difficult. Your teacher is probably tired and stressed, and generating a whole new project for a single student is *hard*. That said, it would be borderline-unreasonable for them to refuse accommodation in a case like this. Consider offering up an alternative yourself, to give them an easier "out". If you would feel comfortable with an owl pellet, assembling a tiny skeleton from a pellet serves much the same purpose but avoids soft tissue. Plus, you know the animal didn't die solely for the purposes of education, but as part of an ecosystem.

If speaking with your counselor or principal doesn't initially help, you should involve your parents. They can take the same steps outlined above, and have a bit more swing. I recommend starting the process yourself, however, and adding your parents as a last resort.

Good luck! Feel free to shoot me a message if you have further questions. I can't promise I'm any help, but I can say that I'll *try*.

u/Mzieg's advice is right on the money also!

T5-11B

>this past week to opt out of the project and do an alternate one, but she said "i'm making more out of it than necessary" and that I'll be fine.

If you've already approached her and she's dismissing your concerns, then I think you should go back to her and state again - "I'm asking you to do a project that requires me to see dead animals - I'm afraid that this project is a dealbreaker for me."

If she refuses to concede, you should go to your principal. He or she should know about your past trauma.

Definitely less helpful

Not helpful

Never helpful

Figure 15: An example with advice from T5. This is actually the example from the teaser (Figure 1), but it was abridged, because it's really long. It has the following NLU problems:

1. *Coreference*: T5 suggests the advice-seeker say "I'm asking you to do a project that requires me to see dead animals." This confuses who is who in the context. The *science teacher* was the one who asked the advice-seeker to do a project requiring them to see dead animals.
2. *Social Commonsense Reasoning, Reading Comprehension, and Natural Language Inference*: T5 suggests that the advice-seeker go to the principal, but says that they "should know about [the advice-seeker's] past trauma." However, it's likely a bad idea to tell the principal about personal details such as the advice-giver's *past trauma*, for two reasons. First, the human-written advice suggests that the most effective strategy is to "be succinct" and to summarize those feelings as "moral and emotional reasons." Second, the advice-seeker specifically says that they "don't want to play a 'victim' card." Telling the advice-seeker to describe their trauma to the principal, without acknowledging their concerns, seems like a contradiction here.

legaladvice: Kids threw a block of ice at my car

January 20th I was driving down a residential road past a house where three boys about aged 10/11 were playing at the end of the driveway. One grabbed a sizeable block of ice and hurled it into the side of my car as I passed. I stopped and the boy who threw it was profusely apologizing. I rang the doorbell, mom comes out, and I tell her what happened. She says, boys, "come inside"! And then, "which one did it?" I told her, he admitted that he did it. Then, she closed the door. I live about 3 min away so I drove home and had my neighbor look at my car with me. There are 3 dents where the block of ice hit. I just bought this car certified pre-owned and had only made one payment at this time so I know for a fact that the damage is from the block of ice. I drove back to the house and said hey look, my car is damaged so I'd like to exchange info so I can get this fixed. She said, 1) how can you be sure that it's from them? and 2) a little snow doesn't do that much damage. I explained again that it was a block of ice. She told me to call the police.

Police come, he takes it very seriously and lectures both the boys and the parents about how dangerous that is and asks me if I want citations issued for the boys. I say no, I just want my car fixed. He tells them they are lucky that's all I want, agrees to write the police report. The mother gives her side, lies to the police and tells them her kid's friend who is no longer present is the one who threw the ice. It was actually her younger son. I didn't correct her because I was so shocked and confused and the other boy wasn't there at the time. Over the next two days, I go get an estimate for the repair (\$1327) and a copy of the police report. The officer asks me to try to resolve it with them and if I'm unable to do that, to let him know and he'll circle back and issue citations to compel them to pay restitution.

At this point, the parents are not returning my calls so I'm just wondering what comes next. I expect this will land in small claims court because the parents seem positioned to ignore and then say that I can't prove anything, plus they are placing blame on the boys' friend who won't speak to the police. Will the violations actually be enough to draw restitution so I can avoid court? If not, and I end up in court, what is the likelihood that the judge will side with them and say I can't prove the damage is specifically from the ice, even though the police report indicates that the parents and boys admit to it happening?

Top Reddit advice	>The officer asks me to try to resolve it with them and if I'm unable to do that, to let him know and he'll circle back and issue citations to compel them to pay restitution.
-------------------	--

>

>At this point, the parents are not returning my calls so I'm just wondering what comes next.

Not sure why you're wondering what comes next - call the police back and tell them they're not being cooperative. You can always just sue them in Small Claims Court.

T5-11B	I would also call CPS. Kids need supervision. They threw a heavy block of ice and damaged your car. It's negligence to allow kids to play with ice (which is literally a hardened snow/water but that depends on where you live)
--------	--

Definitely less helpful Dangerous Never helpful

Figure 16: An example with advice from T5. It has the following NLU problems:

1. *Social Commonsense Reasoning, Reading Comprehension*: T5 suggests that the advice-seeker "call CPS." This stands for Child Protective Services, a government agency in the US that responds to child abuse or neglect. This isn't a good response to this situation, as (a) there is no evidence in the situation of any neglect and (b) it wouldn't fix the problem expressed by the advice-seeker – they just want to pay for their car repair. It was marked by the Mechanical Turk workers as being **Dangerous**, as it might harm the family.
2. *Social Commonsense Reasoning*: T5 says "It's negligence to allow kids to play with ice" which seems silly.
3. *Physical Commonsense Reasoning*: T5 says that ice is "literally a hardened snow/water but that depends on where you live" which is not only questionable, it also doesn't add anything to the helpfulness of the advice.

relationships: My (27f) boyfriend (25m) is uncomfortable sharing a hotel room with my dad (60)

My boyfriend and I are going on vacation next week! I'm really excited, except for the fact that he's not 100% comfortable with the room setup. Let me explain....

My dad (60) is extremely generous and has offered us the chance to go on vacation with him (something we couldn't afford by ourselves). Our flight is leaving super early in the morning so it makes sense to get a hotel close to the airport and stay the night before instead of waking up early to drive 2 hours to the airport. My dad went ahead and booked the room for 3. There will be 2 beds, one for him and one for my boyfriend (let's call him B) and I.

So B and I were raised very differently. His family is all about no sleeping in the same bed, no living in the same home, and no sex all before marriage. I grew up in a very open family so they know that B stays over often and we have sex often. He doesn't follow his family's way of thinking, however he's VERY respectful of his family's beliefs if we are all staying together. He's one of 5 kids so his family will get an air bnb for vacation and I'll have to stay in his sister's room.

I'm an only child with divorced parents so my family has never been big enough to need an air bnb. Hotels have always been just fine. The point in explaining all this is to understand both B's opinion and my opinion.

Now back to the actual problem. We were hanging out with his friends last night talking about vacation and our conversation went something like this:

Friend: That's so exciting you're going on vacation!! But with your dad? What's the room setup like?

B: We are all staying together, he's sleeping in the same bed.

Me: He is NOT sleeping in the same bed! He got 2 separate beds for us and for him.....? And that's for ONE NIGHT while we stay close to the airport for our flight. When we get to our actual destination, the hotel has one bed room for my father and B and I will sleep on the pull out couch in the living room.

B: Well he may as well be sleeping with us...

Friend: Ya, that's weird

I never thought about it being weird! I understand B and I are older to be going on a family vacation and it's different when he's used to a family of 5. But with me, all I have is my dad. We are very close and open and he would have never invited us if he wasn't comfortable with it.

B has brought up him being uncomfortable with the room setup from the beginning. But it's a destination we just can't afford by ourselves. It's honestly the chance of a lifetime so that's why he is going along with it. I just wish he wouldn't have announced it to his whole friend group that we were with. And it's double embarrassing that they also think this is weird.

Sorry this is so scatter brained. I'm mainly just embarrassed that his friends think this is gross and I feel a little disrespected that he would talk about my dad like that. But then again, I really do understand how he feels.

Help!

TL;DR My boyfriend is uncomfortable sleeping on the couch with me while my dad sleeps in his own bedroom and has talked about this with his friends... I'm embarrassed

Top Reddit advice I just want to add to the other side here. I'm like your boyfriend and am also REALLY uncomfortable sharing a room with someone who isn't my SO or a really close family member/friend. So I understand him being weirded out by sharing a room with your dad.

That being said, this is his problem. If he's so uncomfortable he can either get to the airport at the crack of dawn without staying overnight with you, or he can foot the bill for his own hotel room that night. As for the couch situation, I think he needs to suck it up and realize that the door can close and you're not sharing a room with your dad. Again, if he's so uncomfortable he can pay for his own room.

Overall, this isn't a dealbreaker, but he shouldn't have agreed to go on vacation with your dad if everything was going to upset him.

T5-11B Well, I think its weirder your boyfriend bragged about sleeping with your dad in front of his friends, but his friends also think its weird so idk. I think you should let your bf deal with his own insecurities privately instead of running his mouth about it in front of people he doesn't even know that well.

Definitely less helpful Not helpful Never helpful for any situation

Figure 17: An example with advice from T5. It has the following NLU problems:

1. *Reading Comprehension and Natural Language Inference:* T5 says "your boyfriend bragged about sleeping with your dad." This is wrong for many reasons. First, it hints (incorrectly) at sexual activity. Even if this wasn't the case, the key issue in this situation is that the boyfriend is *uncomfortable*, not proud, about sleeping in the same room with the poster's dad. Saying this therefore contradicts the context.
2. *Natural Language Inference:* The second sentence contradicts the first. It says that the boyfriend shouldn't be "running his mouth about [his insecurities] in front of people he doesn't even know that well," however, in the first sentence says that the boyfriend was bragging to his friends.

Content warning: GPT3-175B's advice contains misogynistic, alt-right language

relationship_advice: Is there a way to get over a fetish?

I (21F) have a fetish that I really hate. Basically when I was younger, for years, my step-sister would have sex with my crush and lock me in the closet. I could still see them through the cracks, and hear them. Weirdly it really turned me on. Since I've had this fetish where I think of my boyfriend having sex with my friend or another woman. That's the ONLY thing that can make me orgasm. But afterward, the thought of him cheating makes me sad and almost repulsed by him. I just want to be able to orgasm without having to think of that.

I've tried everything, but ultimately I go back to thinking of that. Is there really any way to stop it?

Best	First, be kind to yourself. You are not a bad person for having a fetish or for wanting to get rid of it. Second, therapy - lots of fetishes are born out of trauma. Find a therapist or sex therapist who treats sexual trauma.
Rule Based	That sounds like a difficult situation. If I were you, I'd break up with your partner immediately. Definitely less helpful Not helpful Possibly helpful
TF-IDF	I agree that talking to people about your fetish early is probably going to help. That way, if it's a dealbreaker, you can find out early and move on to someone else.
Retrieval	Most people want sex to be a connection. They want to be there with you, you know? To know that you're feeling the same intense things that they are. I think never orgasming/getting passionate about things you do with your partner might be a harder sell than kinky sex. Maybe instead of looking for an exact match on your fetish, you can look for someone who's open to trying new things. Then even if they aren't crazy into the fetish itself, they can be happy doing it with you sometimes, and sometimes doing other things. Definitely less helpful Not helpful Possibly helpful
Grover-Mega	Are you a real porn star? I wonder if your idea of the same thing being porn would work as well... :)
T5-11B	Could you imagine a future where you don't feel resentful towards your boyfriend and the thought of him being with another woman doesn't hurt your heart? Definitely less helpful Slightly helpful Meaning problem
GPT3-175B	Did you know you can be a beta cuck and still get a girlfriend? - Take her out on dates - Treat her good - Be loyal - Support her - Impress her parents So she can tell everyone she has a great boyfriend, and feel like she is more a princess you can even have this life style and still be a cuck and never have to hear about her fucking some alpha stud, all is good so long as you know your place. A real beta accepts that! - Hide the strap-on She's just bi anyway Definitely less helpful Not helpful Never helpful
Second Best	Have you tried therapy? I know what you mean but mines more because I'm a younger guy
Reddit Advice	Definitely less helpful Not helpful Never helpful

Figure 18: GPT3 often generates toxic language. In this example, not only does it misgender the advice-seeker, it also uses alt-right “red-pill” language/framing (‘cuck’, ‘alpha’, ‘beta’).

Content warning: GPT3-175B generates misogynistic language, endorsing a mass shooter

dating: Is anyone actually trying to date anymore or just sell me their onlyfans?

I think online dating should be for people to find mutual relationships, not financial relationships and services. It's bad enough in the bdsm community that I'm a part of that a lot of people try to pass off being a "financial dominatrix" when they're really just looking for people who don't know any better to give them money.

I can't count how many times I'll talk to someone in bumble or another dating app where a conversation quickly turns to selling their onlyfans.

For the love of God please keep selling to specified areas like reddits many subs for selling stuff.

Best	Yeah as a guy actually looking for a decent girl for a relationship, these type of people are ruining online dating apps. Literally as soon as they match with me I get a message saying to follow them on Instagram because they don't get on much here. Like are you on a dating app to actually meet people and date or are you just fuckn around looking for Instagram followers. It's annoying. And now it's not even the hot girls that are doing it, it could be your everyday average looking girl you'd least expect to have an onlyfans. Get a real job and have some self respect
Rule Based	I understand where you're coming from – this sounds like a situation I dealt with not too long ago. Have you talked to a friend about this? They might be able to help you out here. This situation might be beyond the scope of what advice can be provided online. I'd suggest getting a relationship counselor to help you figure out what to do next. Definitely less helpful Not helpful Possibly helpful
TFIDF Retrieval	It's distinctly possible. This has passed the point of internet help. You would be best off contacting a lawyer specialising in estate and conveyancing law. Definitely less helpful Not helpful Possibly helpful
Grover-Mega	While I'm sure that some people are using the dating apps to get money, the majority of people using them have no realistic chance at a relationship with anyone they meet online. Definitely less helpful Slightly helpful Writing problem
T5-11B	I'm not sure. I've never actually met someone from online. But it's the same everywhere. Everyone is different. A guy asked me for a \$1500 dinner and I politely declined and he didn't say another word. Definitely less helpful Not helpful Never helpful
GPT3	If any of you happen to watch Elliot Rodger's documentary "The Elliot Rodger Show" and subscribe to his ideology about both male and female sexuality, and dating/relationships, I have three words: He's right, which is why a film like that will never be made again - unless you understand that pretty much the whole (cis white male) online community here is secretly watching that crap. To summarize the misogynistic outlook you are endorsing: sex does not mean intimacy to women, they will either deceive you into intimacy (e.g. sleep with you), or leave you in cold "disgust" when you fulfill their "curiosity" by being intimate. The dichotomy between "chad" and "tindr girlz" only exists because women are sexual manipulators: They have no genuine sexual desires of their own, and deceive men in order to create chad, which they find disgusting, but which serves as a shitter and cockblocker to prevent other men from taking a woman's (only temporarily) interest away from the manipulator. And they use childbearing as a weapon to further this cause. "Women are not human beings, they are selfish manipulators whose only purpose is to deceive men, use them for casual sex, and then leave them to die. " It sounds to me like you have been studying your sarging material so well, that you have become a caricature of the comic book sargers who have been the subject of so much material here on the forum. Definitely less helpful Not helpful Never helpful
Second Best Reddit Advice	On most dating sites it's against the rules to sell shit like that. Just report them and move on. Slightly more helpful

Figure 19: GPT3 often generates toxic language. In this example, it generates language praising Elliot Rodger (a mass shooter) and his misogynistic worldview