

# Shaping the Narrative Arc: Information-Theoretic Collaborative Dialogue

Kory Wallace Mathewson<sup>1\*</sup>, Pablo Samuel Castro<sup>2</sup>, Colin Cherry<sup>3</sup>,

George Foster<sup>3</sup>, Marc G. Bellemare<sup>2</sup>

<sup>1</sup>DeepMind, <sup>2</sup>Google Brain, <sup>3</sup>Google Translate

## Abstract

We consider the challenge of designing an artificial agent capable of interacting with humans in collaborative dialogue to produce creative, engaging narratives. Collaborative dialogue is distinct from chit-chat in that it is knowledge building, each utterance provides just enough information to add specificity and reduce ambiguity without limiting the conversation. We use concepts from information theory to define a *narrative arc* function which models dialogue progression. We demonstrate that this function can be used to modulate a generative conversation model and make it produce more interesting dialogues, compared to baseline outputs. We focus on two modes of modulation: *reveal* and *conceal*. Empirically, we show how the narrative arc function can model existing dialogues and shape conversation models towards either mode. We conclude with quantitative evidence suggesting that these modulated models provide interesting and engaging dialogue partners for improvisational theatre performers.

## Introduction

Designing and building computational models that generate meaningful dialogue for human-interaction is a challenging open problem. Conversational agents can be effective for health-care (Bickmore and Giorgino 2006), by supporting cognitive-behavioural therapy for treating depression (Fitzpatrick, Darcy, and Vierhile 2017), and supporting reminiscence (Nikitina, Callaioli, and Baez 2018) if they are capable of interaction and collaboration.

Rule-based conversational models have existed for over 50 years (Weizenbaum 1966). These methods are limited by hand-tuning and engineering to predict and handle possible inputs. Conversely, generative language models maximize the likelihood of an utterance (e.g. a sentence or sequence of words) (Graves 2013). These models can predict the likelihood of an utterance by considering sentences as a sequences of words or tokens. This objective generates grammatically correct and semantically related to surrounding context it, but lack global consistency (Liu and others 2018).

What makes some dialogues more interesting than others? Interesting collaborative dialogue constructs knowledge iteratively (Swain 2000) and depends on each speaker bringing

information to the conversation (Sawyer 2003). Interest-  
ingness is subjective and difficult to directly optimize via numerical methods (Li and others 2016a).

Our work uses a narrative arc to incrementally construct shared knowledge. A narrative arc defines evolving qualities of emotion, tension, or topic over a story. We draw inspiration from improvised theatre, where actors collaborate in real time to develop narrative based on thematic constraints (Johnstone 1979). Improvised theatre is a unique storytelling medium which relies on collaborative dialogue in which each utterance contributes significant information (Swain 2000). We appeal to the two golden rules of improvised dialogue, characteristic of interesting collaborative dialogue (Johnstone 1979; Sawyer 2003). Good dialogue should 1) accept (i.e. be consistent with the dialogue thus far and 2) reveal (i.e. progress the dialogue with new information).

In this work, we propose a new method to modulate a conversation model, which accepts input utterances by generating consistent and revealing responses. Our approach combines a conversational model with a topic classifier, which we call a universe model. We borrow the term universe from improvised theatre where it is used to describe the world-as-we-know-it (Johnstone 1979; Raby 2010). A universe encompasses associations in the dramatic world and is motivated by the possible world semantics theory (Kripke 1963).

We identify two modes of operation for our shaping method: *revealing* and *concealing*. Revealing dialogue adds additional information about the current universe. Generating utterances which progress a scene with new information is the primary goal of our approach. Concealing dialogue avoids exposing new information about the universe. The ability to generate both revealing and concealing dialogue is a convenient side-effect of this method.

The universe model characterizes the information revealed by each utterance in a sequence. We refer to this information profile across utterances as the narrative arc. By tuning how revealing the model is, we selectively choose utterances to shape the narrative arc to produce more interesting and engaging dialogue. We argue that a balance between revealing and concealing is required for interesting and engaging collaborative dialogue; both over-specification and ambiguity are undesirable. We hypothesize that there is an ideal region of information revelation which our method can expose in existing text-based narratives such as movie scripts.

---

\*This work was done during an internship at Google Brain. Correspondence to: korymath@google.com

## Shaping the Narrative Arc

In this section, we present a mechanism for shaping the narrative arc inspired by combining methods exploring entropy in textual documents (Shannon 1951) with the *Simple Shapes of Stories* described by Vonnegut.<sup>1</sup> We describe concepts of conversation and universe models. Then, we show how these combine to describe a narrative arc. Finally, we show how the narrative arc can be used to generate interesting dialogue.

### The Conversation Model

A conversation model accepts an input utterance and generates one, or several, output utterance(s). The model maintains local coherence by conditioning output generation on the input. We write  $\mathcal{X}$  to denote the set of possible utterances (i.e. sequences of words); in this work,  $\mathcal{X}$  is a collection of English sentences. A sequence of  $t$  successive utterances is a dialogue, denoted  $x_{1:t}$ . A conversation model yields probability  $q$  of utterance  $x_t$  given dialogue  $x_{1:t-1}$ .

We focus on dialogue generation using three retrieval-based conversation models. The first two models are based on the OpenSubtitles dataset (Lison, Tiedemann, and Kouylekov 2018). When queried with an input line  $x_{t-1}$ , a model returns  $K$  candidate responses:

- **Baseline Random model:** sample  $K$  lines from  $\mathcal{X}$ .
- **Deep neural network model (DNN):** we embed all the lines in  $\mathcal{X}$  into a latent semantic space  $S$  using the Universal Sentence Encoder (Cer and others 2018). We encode the input line into  $S$ , and return the  $K$  approximate nearest neighbours in  $S$  using the  $L^2$  norm as the distance metric. Similar to the DNN model, a third model (**Books**), responds with semantically related nearest neighbour lines from literature, filtered for offensive content.<sup>2</sup>

### The Universe Model

The universe model measures how each successive utterance of a dialogue influences the probability distribution over universes. For a given utterance, the universe model calculates a probability distribution over universes. For a sequence of utterances, we use recursive universe belief propagation to update the posterior over the course of a dialogue. Revealing dialogue would concentrate probability mass on a single universe, and concealing dialogue would maintain posterior likelihood over a set of universes. The shape of this sequence of posteriors is the narrative arc. We investigate reveal and conceal dynamics using three different universe models based on probabilistic topic classifiers.

- **Newsgroups:** Using the newsgroup classification dataset, we filter out stop-words, create frequency vectors, and use the TF-IDF (term frequency / inverse document frequency) word weighting scheme to account for word importance in the corpus. We train a naïve Bayes classifier on 5 aggregate topic universes (COMPUTERS, RECREATION, RELIGION, SCIENCE, and TALK).
- **Movies:** naïve Bayes classifier, trained similar to Newsgroups, using a collected dataset of film synopses and one

of 10 corresponding genres (DRAMA, COMEDY, HORROR, ACTION, CRIME, ROMANTIC COMEDY, ROMANCE, THRILLER, FILM ADAPTATION and SILENT FILM) from Wikipedia data (Hoang 2018).

- **DeepMoji:** Deep neural network that takes input text and outputs a distribution over a set of 8 aggregated emoji universes: (SAD, MAD, MEH, NERVOUS, GLAD, MUSIC, LOVE, and MISCELLANEOUS) (Felbo and others 2017). Input text is not transformed, and a pretrained model is used.<sup>3</sup>

### Recursive Universe Belief Propagation

We desire a means by which we can update the universe belief incrementally as evidence is accumulated with each successive utterance in a dialogue. We begin by defining the notion of a *universe model* as a means of modelling the dynamics of information revelation. Consider a finite set of universes,  $\mathcal{U}$ . The role of a universe model is to assess the compatibility of an utterance with a given universe,  $u \in \mathcal{U}$ . Given such a model, we develop a method to update the agent’s posterior universe distribution over a sequence of utterances. For each universe  $u$ , the universe model assigns a likelihood  $p(x_t | x_{1:t-1}, u)$  to an utterance  $x_t$ , conditioned on a dialogue  $x_{1:t-1}$ .

The universe model iteratively updates a posterior belief over universes in a similar spirit to prediction with expert forecasters (Cesa-Bianchi and Lugosi 2006). The probability of a given universe depends on iteratively combining evidence in support of that universe. The posterior probability over universes  $\mathcal{U}$  given a sequence of  $t$  utterances  $x_{1:t}$  is recursively defined as:

$$p_t(u | x_{1:t}) = p_{t-1}(u | x_{1:t-1}) \times \frac{p(x_t | x_{1:t-1}, u)}{p(x_t | x_{1:t-1})}$$

Where  $p_{t-1}(u | x_{1:t-1})$  is prior probability,  $p(x_t | x_{1:t-1}, u)$  is likelihood of utterance conditioned on the past dialogue and universe, and  $p(x_t | x_{1:t-1})$  is likelihood of utterance under the conversation model.

Let  $p_0(u | \cdot) = 1/|\mathcal{U}|$ ,  $u \in \mathcal{U}$  be an initially uniform distribution over universes (i.e. universe model’s prior). We can marginalize out the universe if the evidence is consistent over all hypotheses. To illustrate the relationship between utterance likelihood and universe, we can explicitly write the marginal likelihood as:

$$p(x_t | x_{1:t-1}) = \sum_{u'} p_{t-1}(u' | x_{1:t-1}) p(x_t | x_{1:t-1}, u')$$

Thus, the posterior is updated recursively as:

$$p_t(u | x_{1:t}) = p_{t-1}(u | x_{1:t-1}) \times \frac{p(x_t | x_{1:t-1}, u)}{\sum_{u'} p_{t-1}(u' | x_{1:t-1}) p(x_t | x_{1:t-1}, u')} \quad (1)$$

In practice, it may be convenient to use the output  $z(u | x_t)$  of a probabilistic classifier in lieu of a likelihood function conditioned on past utterances  $x_{1:t}$  and universe  $u$ . Universe classifiers can be trained separately from language models,

<sup>1</sup>From K. Vonnegut lecture: <https://goo.gl/JuEDVR>

<sup>2</sup><https://books.google.com/talktobooks/>

<sup>3</sup>[github.com/bfelbo/DeepMoji](https://github.com/bfelbo/DeepMoji)

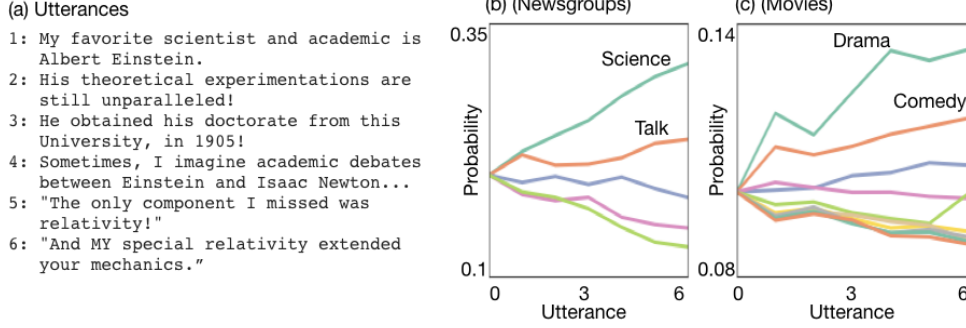


Figure 1: Narrative arcs of synthetic dialogue (a) using the Newsgroups universe model (b) and Movies universe model (c). Dialogue is likely SCIENCE or TALK under the Newsgroups model, and DRAMA or COMEDY under Movies.

and provide complementary signal if model input distributions overlap. This assumption is justified when both models work with similar training corpus vocabularies. We view the probability distribution over universes output by the universe model as derived from a joint distribution  $z(u, x_t)$ , of the universe  $u$ , and utterance  $x_t$ . With  $z(u)$  as the prior distribution over universes, the conditional probability is:

$$z(u | x_t) = \frac{z(u, x_t)}{z(x_t)} = z(u) \times \frac{z(x_t | u)}{z(x_t)}$$

We can substitute  $z(\cdot | x_t)$  for  $p(x_t | x_{1:t-1}, \cdot)$  in Eq. 1 by assuming conditional independence (i.e.,  $p(x_t | x_{1:t-1}, u) = p(x_t | u)$ ), uniform prior distribution (i.e.,  $z(u) = 1/|\mathcal{U}|$ ,  $u \in \mathcal{U}$ ) and constant marginal probability (i.e.,  $z(x_t) = \sum_{u'} p_t(u')p(x_t | u')$ ). These assumptions are justified when the probabilistic topic classifier is a naïve Bayes classifier with uniform prior (Bishop 2006). Thus, the substitution follows the following steps:

$$\begin{aligned} p(x_t | x_{1:t-1}, u) &\approx z(x_t | u) && \text{[cond. independence]} \\ &= \frac{z(u | x_t) z(x_t)}{z(u)} && \text{[Bayes' theorem]} \\ &\approx z(u | x_t) z(x_t) && \text{[}z(u) \text{ uniform prior]} \\ &\approx z(u | x_t) && \text{[}z(x_t) \text{ const. marginal]} \end{aligned}$$

Eq. 1 thus becomes:

$$p_t(u | x_{1:t}) = p_{t-1}(u | x_{1:t-1}) \times \frac{z(u | x_t)}{\sum_{u'} p_{t-1}(u' | x_{1:t-1}) z(u' | x_t)} \quad (2)$$

### The Narrative Arc

As defined in Eq. 2, the posterior  $p_t(\cdot)$  is a function of the dialogue  $x_{1:t}$ . We define the *narrative arc* as the sequence of universe distributions  $p_0(\cdot), p_1(\cdot), \dots$  iteratively calculated for the dialogue. The arc depicts the evolution of a belief over a set of universes. The narrative arc function maps  $\mathcal{X}^t \rightarrow \mathcal{S}(\mathcal{U})^t$ , where  $\mathcal{S}(\mathcal{U})$  is a probability simplex over  $\mathcal{U}$ . We discuss three properties of the narrative arc of the synthetic dialogue shown in Fig. 1:

**1. Utterances affect the arc in varying degrees.** “My favourite scientist and academic is Albert Einstein” is similarly likely under SCIENCE and TALK, and less likely under

the RECREATION universe (bottom green line). Different utterances should have different effects on  $p_t(\cdot)$ .

**2. A concentrating posterior signals a revealing dialogue.** A dialogue which emphasizes scientific content, for example, should see  $p_t(\text{SCIENCE} | \cdot) \rightarrow 1$ . Conversely, we would expect a concealing dialogue to spread the posterior across multiple universes.

**3. A universe model is a perspective on dialogue.** Different universe models can expose different aspects of the same dialogue. Replacing the Newsgroups universe model by a Movies universe model suggests the dialogue is from a DRAMA and/or COMEDY universe. This dialogue would be considered revealing under both universe models.

The universe model can be used to analyze preexisting dialogue, but the model also provides a criterion for favouring utterances when generating dialogue.

### Generating Dialogue with the Narrative Arc

The entropy of the posterior  $p_t(\cdot)$  is given by:

$$H(p_t(\cdot)) := - \sum_{u \in \mathcal{U}} p_t(u) \log p_t(u)$$

Then, the entropy change  $\Delta(\cdot)$  due to a new utterance,  $x_t$ , given the past dialogue,  $x_{1:t-1}$ , is defined as:

$$\Delta(x_t; x_{1:t-1}) := H(p_{t-1}(\cdot)) - H(p_t(\cdot))$$

The term  $\Delta(x_t; x_{1:t-1})$  measures how much a given utterance  $x_t$  changes the entropy of the posterior, given the previous utterances  $x_{1:t-1}$ . A positive value of  $\Delta(\cdot)$  is a reduction in entropy (i.e. information about the universe is revealed). Conversely, a negative value of  $\Delta(\cdot)$  is an increase in entropy (i.e. concealing). We define the score of an utterance  $x_t$ , with respect to a dialogue,  $x_{1:t-1}$ , as:

$$\sigma(x_t; x_{1:t-1}) := \exp\{\alpha \Delta(x_t; x_{1:t-1})\}, \quad \alpha \in \mathbb{R}$$

The exponential function is a convenient way to ensure strict positivity and preserve the ordering of scored candidates. We use our entropy-based score function  $\sigma$  to modulate the sampling of a base conversation model,  $q$ , toward  $\tilde{q}$ , which depends on the change in entropy due to the new utterance.

$$\tilde{q}(x_t | x_{1:t-1}) \propto q(x_t | x_{1:t-1}) \times \sigma(x_t; x_{1:t-1}) \quad (3)$$

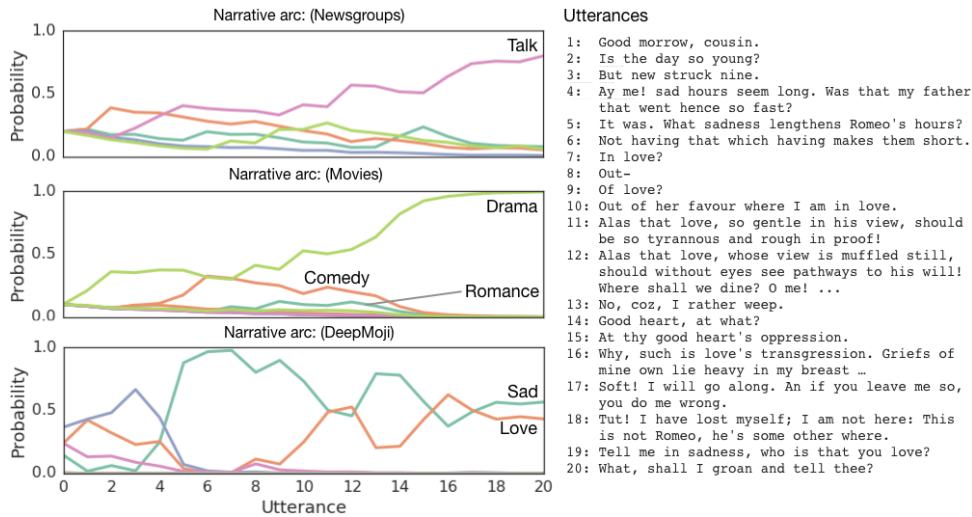


Figure 2: First 20 lines of Romeo and Juliet modelled with NewsGroups (top), Movies (middle), and DeepMojji (bottom) models.

If  $\alpha = 0$ ,  $\sigma(\cdot) = 1$  and candidates are sampled according to  $\tilde{q} = q$ . If  $\alpha \neq 0$ ,  $q$  is modulated by the score  $\sigma(\cdot)$ . Modulation mode depends on the value of  $\alpha$ :

- $\alpha > 0$  (reveal): modulate  $q$  towards revealing the universe. The probability of utterances likely under the universe with highest probability are increased.
- $\alpha < 0$  (conceal): modulate  $q$  towards concealing the universe. The probability of utterances likely under multiple unlikely universes is increased. Utterances not supporting the likely universe are made more likely.

We use these two modulations for filtering samples from our base conversation model. We filter via one of two methods for sampling from an unnormalized distribution: **greedy sampling** and **rejection sampling**. Greedy sampling scores a set of samples from the conversation model and selects the candidate with the maximum score. Scoring a large set of candidates can be time intensive. Rejection sampling (Alg. 1) can sample from the desired unknown modulated distribution online (Murphy 2012). As the entropy function is bounded, the utterance score  $\sigma$  is bounded. In practice, we set a max score and weigh all utterance scores  $\sigma$  above the threshold equally. Both filtering methods have benefits. Rejection sampling provides a smoother distribution and does not require scoring a large set of candidates. Greedy sampling is less sensitive to the range of  $\Delta$  from different utterances.

## Evaluation

### Narrative Arc of Existing Dialogues

In Fig. 2, we visualize the narrative arc underlying the first 20 lines of Shakespeare’s Romeo and Juliet using three universe models: 1) NewsGroups, 2) Movies, and 3) DeepMojji.

Fig. 2 illustrates the entropy-reducing nature of good dialogue by showing us the underlying, evolving, narrative arc. Under the NewsGroups universe model, the dialogue evolves toward a TALK-centric universe. Under the Movies model, the same dialogue balances between comedy and drama before shifting towards drama. Finally, using the DeepMojji

---

### Algorithm 1 Generating dialogue with rejection sampling.

---

**Given:** conversation model  $q$ , scoring function  $\sigma$ , first line  $x_1$ , length  $N$ , max score  $M$ , max samples  $S$

**Return:** dialogue  $x_{1:N}$

**for**  $t$  in  $2 \dots N$  **do**

**while**  $\text{step} \leq S$  **do**

sample  $x_t \sim q(x_t | x_{1:t-1})$

sample  $r \sim \text{Uniform}(0, 1)$

**if**  $r \leq \sigma(x_t; x_{1:t-1})/M$  **then**

append  $x_t$  to  $x_{1:t-1}$

**break**

**end if**

**end while**

**end for**

---

universe model, a developing ambiguity between DeepMojji universes SADNESS and LOVE is uncovered. This supports the hypothesis that existing dialogues exhibit underlying narrative arcs conditioned on universe models.

### Shaping the Narrative Arc

In this section, we demonstrate that our method is able to modulate conversation models toward generation of revealing or concealing dialogues. Linguistic quality and semantic consistency of utterances are determined by the language underlying the conversation model. We emphasize evaluation of narrative arc shaping by focusing on the information contribution of the subsequent utterances.

We use the DNN conversation model to test how preferential selection, induced by our score function, can modulate information introduced into the conversation. In Fig. 3 we present characteristic narrative arcs and dialogues using concealing (top), neutral (middle), and revealing (bottom) modes. Each generation was primed with the first two lines from Romeo and Juliet (shown in bold in Fig. 3).

A significant difference is exposed between concealing (top) which tends toward a high entropy, uniform universe

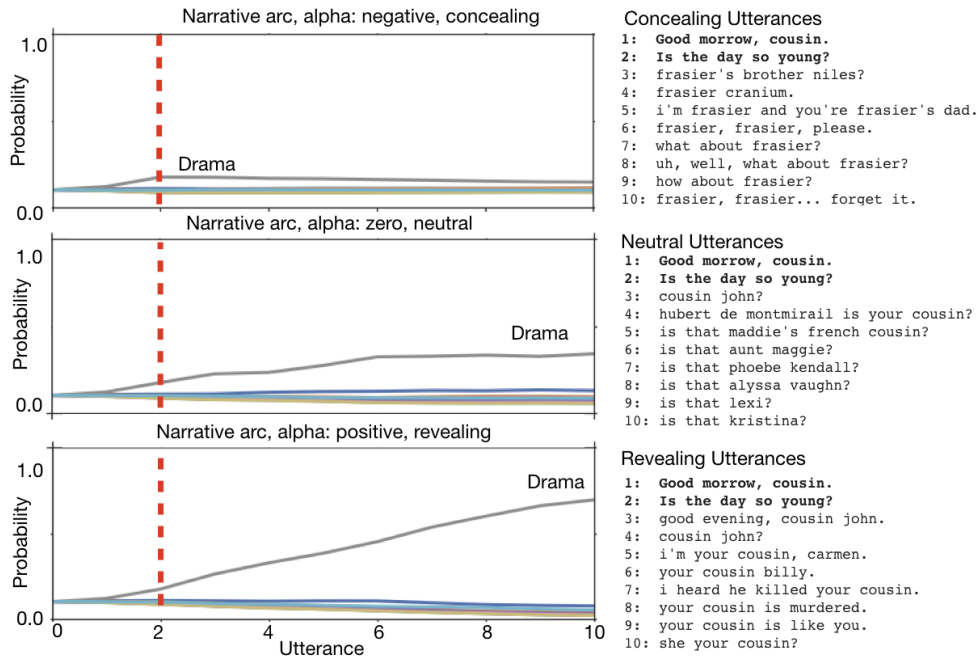


Figure 3: Narrative arcs over 10 utterances at increasing  $\alpha$  values: concealing (top), neutral (mid), revealing (bottom). On the right are utterances generated by each model after priming (bold). Dotted red line indicates the start of narrative arc shaping.

distribution, and revealing (bottom) where drama tends toward 1.0. DRAMA remains the most likely universe (and visible on all plots) as it was supported by the first two lines and subsequent utterances did not significantly shift the distribution. Fig. 3 also shows the utterances selected by the model. Concealing utterances do not add information to the dialogue, revealing utterances incorporate new information over the course of the dialogue.

We next evaluate our method’s ability to filter for concealing/revealing utterances by measuring the entropy under both an objective universe (i.e. the universe model used for scoring in generation) and a test universe not used for scoring. We use the Newsgroups universe model for objective scoring and the Movies model for testing. A random conversation model is used to generate response candidates.

We generate 20 conversations following a process similar to Algorithm 1 but using greedy sampling. Each conversation starts with a random dialogue starter line to encourage diversity and then 19 lines are sampled from the conversation model using the narrative arc function. This approximates the length of a medium-duration improvised conversation (Sawyer 2003).

Results are presented in Fig. 4. There is a significant difference between the entropy under the objective and testing universes, but each model exhibits similar dynamics over the dialogues. We conclude that concealing dialogue can conceal under multiple universes, and revealing dialogue can reveal information under multiple universe models.

The revealing/concealing dynamics of each utterance may be related to measurable lexicographical qualities such as words per sentence (WPS). We analysed the language used in 190 lines from each model and found a significant difference

( $p < 0.001$ ) between utterances selected by the revealing model ( $9.26 \pm 5.7$  WPS) and utterances selected by the concealing model ( $5.05 \pm 2.79$  WPS).

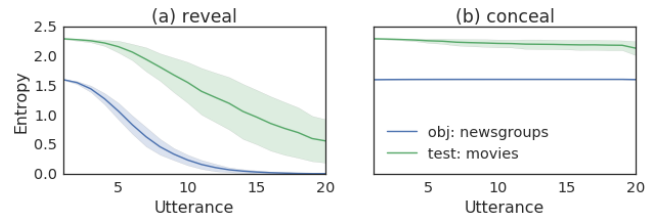


Figure 4: Revealing and Concealing across Universe Models. Dialogue generated to be (a) revealing ( $\alpha = 20$ ) under the objective model Newsgroups is revealing under the testing Movies universe. The same is true for (b) concealing ( $\alpha = -25$ ) dialogue. Data shown are means and standard deviation (shaded) over 20 runs of random conversation model.

### Predicting the Next Best Line

We next test the system’s ability to add information to improve performance on a prediction task. Given a sequence of 5 gold-standard conversational utterances and a list of 10 next utterance candidates (i.e. the ground truth and 9 distractors), can the universe model be used to improve accuracy of predicting the ground truth?

Evaluation compares top-3 accuracy and mean reciprocal rank (MRR) over samples in a held out test set. Accuracy measures the likelihood that the system scores the ground truth within the top-3 candidates against the distractors. MRR

compares average ground truth ranking across conditions. A Transformer language model was trained on OpenSubtitles (Lison, Tiedemann, and Kouylekov 2018) to predict an output given a set of input lines (Vaswani and others 2017).

The trained Transformer model was used to assign a perplexity score for output line candidates given an input context line. For each unique subtitle file in the validation and test sets, the concatenation of the first 5 lines serve as input context and line 6 is the ground truth output to be predicted. Negative candidates are randomly selected from lines in the respective corresponding data segment (i.e. validation/test), thus may not be from the same file as the input context lines.

The perplexity under the trained conversation model serves as the unmodulated probability  $q(x_t|x_{1:t-1})$  (Eq. 3) of selection in the prediction task. The input sequence is then passed, line-by-line, through a Newsgroups universe model and a score is assigned to each candidate relative to the change in entropy of the evolving posterior. The  $\alpha$  value is modulated over 100 evenly spaced values between  $[-2, 2]$ . The accuracy of predicting the ground truth in the top-3 candidates and the MRR of the ground truth are computed.

The results on the validation set are shown in Fig. 5. By selecting the correct  $\alpha$  value, the likelihood of correctly selecting utterances revealing an incremental amount of information increases significantly. Note the shape of the curve as  $\alpha$  changes. As hypothesized, there exists a region, between 0 and 1 where the ‘right’ amount of universe information is revealed. This region corresponds to the notion that each line of dialogue will reveal some, but not too much, information about the universe. As  $\alpha$  continues to increase, the accuracy decreases below the neutral baseline. The top-3 accuracy of prediction increases when the universe model boosts the probabilities of appropriately revealing dialogue. The validation set is used to set the optimal  $\alpha$ , which is then used to score samples in the test set and results are presented in Table 1. Two additional models are included for comparison. *T2T@1* uses 1 preceding the ground truth as context. *Unigram* assigns a perplexity to output candidates by building a unigram language model using the 5 input lines as a corpus. A smoothing factor of  $1 \times 10^{-9}$  is used for out-of-vocabulary words. Additionally, a random conversation baseline model is included. For each model tested, information from the universe model significantly improves the predictive accuracy on this task.

### Interactive Collaborative Dialogue

Finally, as a practical implementation case-study, we tested how this system performs in collaborative dialogue through interaction with humans. 4 expert improvisational theatre performers engaged with the system in 3 text-based conversations. Each conversation consisted of 5 utterance-response pairs for a total of ten utterances (i.e. an average length of a short-duration improvised scene (Sawyer 2003)). Subjects are native English speakers with 5+ years professional performance experience and are familiar with shared narrative development and collaborative dialogue. Each interacted with revealing, concealing, and neutral models in a randomized order unknown to the them.

This experiment used the Books conversation model and

CM	UM	Top3Acc	MRR
T2T@5	NG	<b>0.520</b>	<b>0.456*</b>
T2T@5	Neutral	0.507	0.444
T2T@1	NG	0.483	0.428*
T2T@1	Neutral	0.469	0.412
Unigram	NG	0.366	0.337*
Unigram	Neutral	0.296	0.290
Random	Neutral	0.302	0.294

Table 1: Results for predicting the next line. CM is the conversation model, UM is the universe model, Top3Acc is the accuracy of predicting the ground-truth in the top-3 of 10 candidates, and MRR is the mean reciprocal rank of the ground truth. Unigram CM calculates the perplexity of each candidate given the input lines as training corpus. T2T@N is a Tensor2Tensor Transformer model which uses the previous N lines as an input to predict the output and NG is the Newsgroups universe. A Neutral universe model represents no modulation which is equivalent to  $\alpha = 0$ . \* indicates  $p < 0.05$  for a Students’ t-test comparing MRR to the Neutral model.

the DeepMoji universe model. Following the interactions, each performer was asked the following question: “please rank the conversations from 1 (most engaging) to 3 (least engaging)”. Engagingness was defined to align with the notions of revealing and concealing in this work. An agent is engaging for shared scene development if it brings just enough information to add specificity and reduce ambiguity but not limit the conversation.

Three of the four performers ranked the revealing model,  $\alpha > 0$ , as the most engaging. Those three performers ranked  $\alpha = 0$  as being less engaging due to being “too random”. All subjects ranked  $\alpha < 0$  as being least engaging and not bringing enough information to the scene. These results support the hypothesis that  $\alpha$  effectively modulates collaborative dialogue engagingness in human-machine interaction.

### Related Work

Collaborative dialogue between humans and machines has been proposed as a grand challenge in artificial intelligence (Mathewson and Mirowski 2017a; Martin, Harrison, and Riedl 2016; Brown 2008). Previous methods have used hard coded rules, decision trees, and event representations to generate novel narrative chains (Martin and others 2017). We use a deep neural network-based generative language model enhanced with universe model information in the context of improvised theatre (Mathewson and Mirowski 2017b).

While neural response generation systems provide a trainable end-to-end system for language generation, these methods are prone to providing generic, unspecific responses (Li and others 2015). Advances have improved generated responses by optimizing sentence encoding and decoding jointly, post-generation candidate re-scoring (Bordes, Boureau, and Weston 2016; Vinyals and Le 2015; Sordani and others 2015), reinforcement learning (Li and others 2016a), hierarchical models for distilling extended context (Serban and others 2016), and auxiliary training objectives, such as maximizing mutual information (Li and others 2015), and personality specificity and consistency (Li

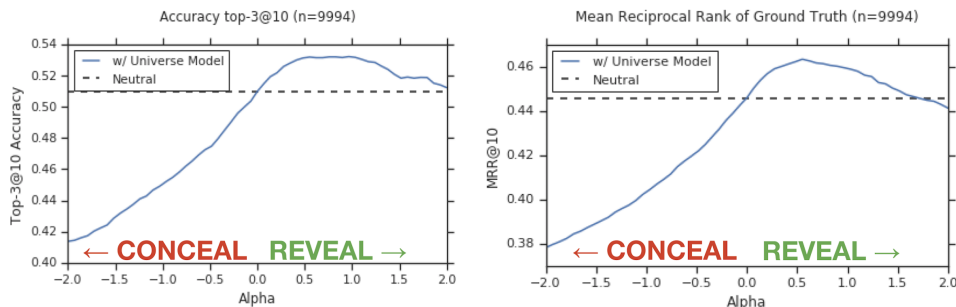


Figure 5: Varying  $\alpha$  for (left) top-3 accuracy, (right) mean reciprocal rank in universe model modulated prediction task

and others 2016b; Zhang and others 2018). In future work, universe models and conversational models could be trained jointly.

Our work is related to the controlled generation of text using disentangled latent representations (Hu and others 2017). Previous work has used a topic-transition generative adversarial network to enforce smoothness of transition of subsequent utterances (Liang and others 2017). These methods use neural encoder-decoders and generate responses given an input sequence and a desired target class for the response.

Other work has aimed to improve candidates returned by retrieval-based conversation models (Weston, Dinan, and Miller 2018). These methods utilize a conversation model to find similar prototypes using embedding distances and refine prototypes with a sequence-to-sequence model (Guu and others 2017). We do not refine candidates from the conversation model, rather we sample and select using a scoring function defined by the revealing and concealing parameter.

Similar to universe models, topic models or lexical fields have been shown capable of tracking general subjects of a text (Blei, Ng, and Jordan 2003). Dynamic topic models characterize the evolution of topics over a set of documents over time (Blei and Lafferty 2006). Our work differs in that we generate dialogue using the evolving probabilistic belief during a single conversation, as opposed to tracking topical shifts over longer time-scales. Using a probabilistic classifier for narrative tracking has been explored previously (Mohammad 2011; Reagan and others 2016). These works used sentiment classifiers to track emotion and plots arcs through narratives. We extend these works by using probabilistic universe models collaborative dialogue generation. While our work uses separate language and universe models, ongoing research aims to steer or control the properties of text generated with language models (Radford and others 2019) using various attribute models during training (Dathathri and others 2019; Saleh and others 2019).

## Discussion and Summary

While innovations have improved the linguistic quality, semantic alignment, and consistency of utterances generated by neural models, generated conversations still lack interestingness and engagingness. Our work selects engaging utterances by shaping the underlying narrative arc as opposed

to improving the training of generative language models. The methods presented are agnostic to both the universe and the conversational model used. Using rules from improvised theatre, we quantitatively define the evolution of interesting and engaging dialogue.

In this work we focus on genre, emoji, and topic-based universe models. Other universe models to be explored involve causality of events, directions of relationships, or audience reaction prediction. While this work explores the interaction between a base conversation model and a universe model, this method could be compatible with image or video generation.

The main contribution of this work is the computational formalization of the narrative arc, an information-theoretic framework for collaborative dialogue interaction. The framework fills a gap in previous research by connecting the utterance-level improvements of language models with the conversation-level improvements of universe tracking. This is done by sampling candidates from a conversational model using a universe model and the narrative arc. We illustrate narrative arcs underlying popular dialogues and show how universe models can be combined with conversation models to aid in interesting dialogue generation. We present empirical results showing how the narrative arc can improve accuracy on a next line prediction task. Finally, we present an expert user-study to validate our model.

## References

- Bickmore, T., and Giorgino, T. 2006. Health dialog systems for patients and consumers. *Journal of Biomedical Informatics* 39(5):556 – 571. Dialog Systems for Health Communications.
- Bishop, C. M. 2006. *Pattern Recognition and Machine Learning*. Springer, New York.
- Blei, D. M., and Lafferty, J. D. 2006. Dynamic topic models. In *Proceedings of the 23rd international conference on Machine learning*, 113–120. ACM.
- Blei, D. M.; Ng, A. Y.; and Jordan, M. I. 2003. Latent dirichlet allocation. *JMLR* 3(Jan):993–1022.
- Bordes, A.; Boureau, Y.-L.; and Weston, J. 2016. Learning end-to-end goal-oriented dialog. *arXiv preprint arXiv:1605.07683*.
- Brown, K. 2008. The auslander test: or, ‘of bots and humans’.

- International Journal of Performance Arts and Digital Media* 4(2-3):181–188.
- Cer, D., et al. 2018. Universal sentence encoder. *arXiv preprint arXiv:1803.11175*.
- Cesa-Bianchi, N., and Lugosi, G. 2006. *Prediction, Learning, and Games*. New York, NY, USA: Cambridge University Press.
- Dathathri, S., et al. 2019. Plug and play language models: a simple approach to controlled text generation. *arXiv preprint arXiv:1912.02164*.
- Felbo, B., et al. 2017. Using millions of emoji occurrences to learn any-domain representations for detecting sentiment, emotion and sarcasm. *ArXiv e-prints*.
- Fitzpatrick, K. K.; Darcy, A.; and Vierhile, M. 2017. Delivering cognitive behavior therapy to young adults with symptoms of depression and anxiety using a fully automated conversational agent (woebot): A randomized controlled trial. *JMIR Ment Health* 4(2):e19.
- Graves, A. 2013. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*.
- Guu, K., et al. 2017. Generating sentences by editing prototypes. *arXiv preprint arXiv:1709.08878*.
- Hoang, Q. 2018. Predicting movie genres based on plot summaries. *arXiv preprint arXiv:1801.04813*.
- Hu, Z., et al. 2017. Toward controlled generation of text. In *International Conference on Machine Learning*, 1587–1596.
- Johnstone, K. 1979. *Impro. Improvisation and the theatre*. Faber and Faber Ltd.
- Kripke, S. A. 1963. Semantical analysis of modal logic i normal modal propositional calculi. *Mathematical Logic Quarterly* 9(5-6):67–96.
- Li, J., et al. 2015. A diversity-promoting objective function for neural conversation models. *CoRR abs/1510.03055*.
- Li, J., et al. 2016a. Deep Reinforcement Learning for Dialogue Generation. *ArXiv e-prints*.
- Li, J., et al. 2016b. A persona-based neural conversation model. *arXiv preprint arXiv:1603.06155*.
- Liang, X., et al. 2017. Recurrent topic-transition gan for visual paragraph generation. *CoRR, abs/1703.07022* 2.
- Lison, P.; Tiedemann, J.; and Kouylekov, M. 2018. Opensubtitles 2018: Statistical rescoring of sentence alignments in large, noisy parallel corpora. In *Proceedings of the 11th edition of the Language Resources and Evaluation Conference (LREC 2018)*, Miyazaki, Japan.(accepted).
- Liu, P. J., et al. 2018. Generating wikipedia by summarizing long sequences. *arXiv preprint arXiv:1801.10198*.
- Martin, L. J., et al. 2017. Improvisational storytelling agents. In *NeurIPS 2017 Workshop on Machine Learning for Creativity and Design*.
- Martin, L. J.; Harrison, B.; and Riedl, M. O. 2016. Improvisational computational storytelling in open worlds. In *International Conference on Interactive Digital Storytelling*, 73–84. Springer.
- Mathewson, K. W., and Mirowski, P. 2017a. Improvised comedy as a turing test. *CoRR abs/1711.08819*.
- Mathewson, K. W., and Mirowski, P. 2017b. Improvised theatre alongside artificial intelligences. In *AAAI AIIDE*.
- Mohammad, S. 2011. From once upon a time to happily ever after: Tracking emotions in novels and fairy tales. In *Proc. of the 5th ACL-HLT Workshop on Language Technology for Cultural Heritage, Social Sciences, and Humanities*, 105–114. ACL.
- Murphy, K. P. 2012. *Machine Learning: A Probabilistic Perspective*. The MIT Press.
- Nikitina, S.; Callaioli, S.; and Baez, M. 2018. Smart conversational agents for reminiscence. *arXiv preprint arXiv:1804.06550*.
- Raby, G. 2010. Improvisation and devising: The circle of expectation, the invisible hand, and rsvp. *Canadian Theatre Review* 143(1):94–97.
- Radford, A., et al. 2019. Language models are unsupervised multitask learners. *OpenAI*.
- Reagan, A. J., et al. 2016. The emotional arcs of stories are dominated by six basic shapes. *EPJ Data Science* 5(1):31.
- Saleh, A., et al. 2019. Hierarchical reinforcement learning for open-domain dialog. *arXiv preprint arXiv:1909.07547*.
- Sawyer, R. K. 2003. *Improvised dialogues: Emergence and creativity in conversation*. Greenwood Publishing Group.
- Serban, I. V., et al. 2016. Building end-to-end dialogue systems using generative hierarchical neural network models. In *AAAI*, volume 16, 3776–3784.
- Shannon, C. E. 1951. Prediction and entropy of printed english. *Bell Labs Technical Journal* 30(1):50–64.
- Sordoni, A., et al. 2015. A neural network approach to context-sensitive generation of conversational responses. *arXiv preprint arXiv:1506.06714*.
- Swain, M. 2000. The output hypothesis and beyond: Mediating acquisition through collaborative dialogue. *Sociocultural Theory and Second Language Learning* 97:114.
- Vaswani, A., et al. 2017. Attention is all you need. *CoRR abs/1706.03762*.
- Vinyals, O., and Le, Q. 2015. A neural conversational model. *arXiv preprint arXiv:1506.05869*.
- Weizenbaum, J. 1966. Eliza: a computer program for the study of natural language communication between man and machine. *Communications of the ACM* 9(1):36–45.
- Weston, J.; Dinan, E.; and Miller, A. H. 2018. Retrieve and refine: Improved sequence generation models for dialogue. *arXiv preprint arXiv:1808.04776*.
- Zhang, S., et al. 2018. Personalizing dialogue agents: I have a dog, do you have pets too? *CoRR abs/1801.07243*.