GAMMT: Generative Ambiguity Modeling Using Multiple Transformers

Xingcheng Xu*

Abstract

We introduce a novel model called GAMMT (Generative Ambiguity Models using Multiple Transformers) for sequential data that is based on sets of probabilities. Unlike conventional models, our approach acknowledges that the data generation process of a sequence is not deterministic, but rather ambiguous and influenced by a set of probabilities. To capture this ambiguity, GAMMT employs multiple parallel transformers that are linked by a selection mechanism, allowing for the approximation of ambiguous probabilities. The generative nature of our approach also enables multiple representations of input tokens and sequences. While our models have not yet undergone experimental validation, we believe that our model has great potential to achieve high quality and diversity in modeling sequences with uncertain data generation processes.

1 Introduction

Risk and ambiguity are two types of uncertainty that are commonly modeled. Risk refers to situations where the probabilities of various outcomes are known or can be estimated, while ambiguity refers to scenarios where the probability distribution of potential outcomes is unknown or uncertain.

In this paper, we present a novel approach for modeling data generation processes under ambiguity, which we call Generative Ambiguity Models using Multiple Transformers (GAMMT). Traditional machine learning methods assume data is sampled from a deterministic probability, but our proposed approach acknowledges that the process may be ambiguous and determined by multiple probabilities. Our GAMMT model utilizes multiple parallel transformers to model a set of probabilities in sequential data, such as token sequences in natural language. These parallel transformers function as a system to respond to ambiguity and provide multiple representations of the underlying input tokens and sequence.

To illustrate the concept of ambiguity, we use Ellsberg's urns, which were introduced in 1961 [13] and have stimulated research on ambiguity in economics and finance. Our GAMMT model

^{*}Shanghai AI Laboratory. Email: xingcheng.xu18@gmail.com

leverages the effectiveness of transformers, which are attention-based mechanisms in deep neural networks, to model probabilities. Each parallel transformer receives an input sequence and outputs a deterministic probability, and the transformers are connected through a selection mechanism.

Our proposed GAMMT models have two key features that distinguish them from existing approaches. Firstly, they are generative and can effectively model ambiguity in data generation processes. Secondly, the last hidden layer of each parallel transformer provides multiple embeddings of each token and the input sequence, enabling diverse representations of the underlying input. We introduce the conceptual framework of our GAMMT models and anticipate that they will improve the quality and diversity of generated text, making it more engaging and human-like. Although we have not yet provided experimental validation, we believe that our models have great potential in this regard. Our approach to modeling ambiguity in sequential data is novel and provides a more interpretable and flexible method than existing techniques. The unique features of our model, such as the use of multiple transformers and its generative nature, make it a promising direction for future research in machine learning. We hope that this framework will stimulate further research in the fields of machine learning and natural language processing.

The paper is organized as follows: In Section 2, we review related work on transformers and ambiguity modeling. In Section 3, we present our proposed notation for capturing the structure of uncertainty and define the set \mathcal{P}^{LM} to represent the set of plausible probability laws. We then introduce the GAMMT architecture in Section 4, which utilizes multiple parallel transformers linked by a selection mechanism to approximate ambiguous probabilities. In Section 5, we provide an overview of the algorithms for the GAMMT model's architecture, model training, and inference. Finally, in Section 6, we conclude this paper.

2 Related Work

Transformers have emerged as a powerful and versatile tool in machine learning, driving significant progress in natural language processing (NLP), as demonstrated by models such as BERT (Devlin et al., 2018 [10]), OpenAI GPT (Radford et al., 2018 [23]), GPT-2 (Radford et al., 2019 [24]), GPT-3 (Brown et al., 2020 [3]), InstructGPT (Ouyang et al., 2022 [19]), ChatGPT ([20]), and Google T5 (Raffel et al., 2019 [26]), PaLM (Chowdhery et al., 2022 [8]), etc. These models have achieved groundbreaking results and transformed NLP research. Additionally, transformers have extended their reach beyond NLP to other fields such as computer vision, exemplified by the success of Vision Transformers proposed by Dosovitskiy et al. in 2020 [11], and even structural biology, as shown by AlphaFold 2 developed by Jumper et al. in 2021 [17]. This broad application of transformers highlights their versatility and potential for future breakthroughs.

Ambiguity is a widespread challenge in various domains, including vision, language, learning, and decision-making. In machine learning, several studies have explored the concept of ambiguity, which can arise from external sources or incomplete knowledge during the generation process.

Works by Ek et al. (2008) [12], Patel et al. (2019) [21], Yang et al. (2021) [35], Buisson et al. (2022) [4], and others have previously explored ambiguity in machine learning. Ambiguity is pervasive in various fields, such as vision, natural language, general learning, and decision-making.

For example, in images, visual relationships can be semantically ambiguous, as classified by Yang et al. [35] into three types: synonymy ambiguity, hyponymy ambiguity, and multi-view ambiguity. In natural language, perceptual vagueness can occur at different levels of granularity, including words, sentences, paragraphs, and documents. Ambiguity in images or language can arise from external sources or due to individuals' feelings, perception, experiences, or incomplete knowledge during the generation process.

Traditionally, machine learning models have represented images, words, or sentences using a fixed feature vector in a deterministic way. However, such a fixed representation is insufficient to capture the ambiguity mentioned above. To address this limitation, researchers have proposed using multiple complementary representations of an underlying phenomenon. For example, Ek et al. [12] demonstrated the use of multi-modal regression on a benchmark human pose estimation dataset. Yang et al. [35] utilized a probability distribution to represent each union region of an image using Gaussian embedding. In this approach, each union region is parameterized by a mean and variance, where the mean vector represents the normal feature vector, and the variance measures feature uncertainty.

While existing approaches, such as those proposed by Ek et al. [12] and Yang et al. [35], have made progress in handling ambiguity in machine learning, they do not explicitly model the probabilistic structure of ambiguity as in our proposed GAMMT model. Our model utilizes multiple parallel transformers to model probabilities in sequential data, providing a more interpretable and flexible approach to handling ambiguity. Additionally, the GAMMT model is generative, allowing it to effectively characterize ambiguity in the data generation process.

Compared to existing approaches, our proposed GAMMT model has several unique features that make it a promising avenue for future research in the field of machine learning. Firstly, it provides multiple ambiguous representations of the input tokens and sequence, allowing for greater flexibility in downstream tasks such as language understanding and generation. Secondly, the use of multiple transformers allows for a more comprehensive modeling of the complex and multi-dimensional nature of ambiguity in sequential data. Lastly, the generative nature of the GAMMT model provides additional insights into the generation process of ambiguous data.

Furthermore, in the field of neuroscience, cognitive neuroscientists have utilized neuroimaging techniques, such as functional magnetic resonance imaging (fMRI), to investigate brain activation and neuronal correlates in humans during decision-making under ambiguity. Studies conducted by Levy et al. (2010) [18], Bach et al. (2011) [2], Chumbley et al. (2012) [9], Taya (2012) [28], and others have demonstrated that a distributed network of brain regions associated with cognitive control and reward processing is closely linked to decision-making in ambiguous situations. Notably, Chumbley et al. [9] observed clear ambiguity-dependent responses in the hippocampus, suggesting that certain neuronal systems may play a role in resolving ambiguity.

This suggests that our GAMMT model, which is designed to model ambiguity in sequential data using multiple parallel transformers, may be able to simulate the distributed nature observed in cognitive neuroscience studies.

3 Framework

In the field of sequential modeling, the standard approach in machine learning involves factorizing joint probabilities P(x) of a sequence $x = (s_1, s_2, \dots, s_\ell) \in V^*$, where $V^* = \bigcup_{\ell=0}^{\infty} V^{\ell}$ is a sequence space consisting of tokens from a vocabulary V. This factorization is commonly achieved by taking the product of conditional probabilities, as in popular language models such as GPT [23], GPT-2 [24], and GPT-3 [3], etc. The joint probability can be expressed as:

$$P(x) = \prod_{n=1}^{\ell} P(s_n | s_1, s_2, \dots, s_{n-1}), \tag{3.1}$$

where $P(s_n|s_1, s_2, \dots, s_{n-1}) := P(s_1)$ when n = 1.

For the conditional probability P(x|z) of a pair of sequences $x = (s_1, s_2, \dots, s_\ell) \in V^*$ and $z \in V^*$, it can be factorized by taking the product of following conditionals, as in the original Transformer model [32]:

$$P(x|z) = \prod_{n=1}^{\ell} P(s_n|s_1, s_2, \dots, s_{n-1}, z),$$
(3.2)

where $P(s_n|s_1, s_2, \dots, s_{n-1}, z) := P(s_1|z)$ when n = 1.

In this section, we propose a novel approach to address sequential modeling problems. We recognize that ambiguity is prevalent in various scenarios, and the randomness of sequences may not be accurately described by a single probability distribution. Instead, we introduce a collection of probabilities to capture the structure of uncertainty associated with probabilities and develop a novel framework for sequential modeling that takes into account this collection of probabilities. Building on the work of Chen and Epstein (2022) [6], we present notation that captures the structure of uncertainty and define the set \mathcal{P}^{LM} (see Section 3.2) to represent the set of plausible probability laws. Our approach recognizes that different sets of probabilities may be associated with different time points, allowing for a more complex modeling of sequential data.

3.1 Uncertainty of Probabilities

In this subsection, we propose a novel approach to modeling uncertainty in sequential modeling problems by utilizing a collection of probabilities. This is analogous to having a mystical box containing multiple distinct dice, rather than just one. For each event or decision, we randomly select and throw one or more of these dice, and the resulting outcome is based on their combined outcomes. The data generation process is thus determined by a set of probabilities denoted by:

$$\mathcal{L}_n = \{ P_{n1}, P_{n2}, \cdots, P_{nM} \}, \quad n = 1, 2, \cdots$$
 (3.3)

where \mathcal{L}_n represents the set of fundamental probability measures for creating the data process. The number of measures in the set determines the degree of model ambiguity or uncertainty. This novel approach enables us to model ambiguity in a more flexible and interpretable way than existing techniques.

The objective of sequence modeling is to maximize the product of conditional probabilities while also incorporating a selection mechanism, which is given by the following expression:

$$\prod_{n=1}^{\ell} S_n \left(\left\{ P_{nj}(s_n | s_1, s_2, \cdots, s_{n-1}), P_{nj} \in \mathcal{L}_n \right\} \right), \tag{3.4}$$

where $P_{nj}(s_n|s_1,s_2,\cdots,s_{n-1}):=P_{nj}(s_1)$ when n=1. The set \mathcal{L}_n represents the collection of plausible probability laws for each one-step-ahead conditional at time point n. To model the collection of conditional probabilities \mathcal{L}_n , deep neural networks with learnable parameters θ_n are used, such as multiple Transformer decoders for natural language models in our paper.

The selection mechanism S_n is used to address the ambiguity structure present in the data. It may involve selecting probabilities at random subject to a uniform random variable or selecting the maximum of all possible conditionals. Notably, our proposed approach degenerates to the standard approach when \mathcal{L}_n is a singleton set $\{P_n\}$.

Our GAMMT model captures the probabilistic structure of uncertainty by utilizing the set \mathcal{P}^{LM} (in the next subsection), which allows for different sets \mathcal{L}_n of plausible probability laws for each one-step-ahead conditional at every time point. By incorporating this more flexible set of probability laws, the GAMMT model can capture a wider range of uncertainty and account for potential dependencies between experiments. In the next subsection, we will provide a detailed description of the probabilistic structure of uncertainty in our GAMMT model.

3.2 Probabilistic Structure of Ambiguity

In this subsection, we present notation that captures the structure of uncertainty associated with probabilities, building on the work of Chen and Epstein (2022) [6]. When there is ambiguity about the probability law for each position in a sequence, understanding the relationship between positions becomes crucial, as this relationship can have significant implications for the interpretation of the results.

Consider a filtered space $(\Omega, \{\mathcal{G}_n\}_{n=1}^{\infty}))$, where $\Omega = \prod_{i=1}^{\infty} \Omega_i$ and each Ω_i represents the set of possible outcomes for the *i*th position in a sequence. Let $\Omega^{(n)} := \prod_{i=1}^n \Omega_i$ and $\Omega_{(n+1)} := \prod_{i=n+1}^{\infty} \Omega_i$. For each n, the σ -algebra \mathcal{G}_n represents the observable tokens regarding positions $1, 2, \cdots, n$ on the underlying probability space $(\Omega^{(n)}, \mathcal{G}_n)$. We assume that \mathcal{G}_n is increasing with n, with \mathcal{G}_0 being the trivial σ -algebra. The σ -algebra generated by all the observable tokens is denoted by \mathcal{G} and is defined as the smallest σ -algebra that contains the sets \mathcal{G}_n , where n ranges over the non-negative integers. In other words, \mathcal{G} is the σ -algebra generated by the union of all the \mathcal{G}_n 's. The precise ex ante probabilities of observable tokens are unknown and are instead represented by a set \mathcal{P} of probability measures on (Ω, \mathcal{G}) .

We provide the following additional notation and terminology:

$$\boldsymbol{\omega}^{(n)} = (\boldsymbol{\omega}_1, \cdots, \boldsymbol{\omega}_n) \in \Omega^{(n)}, \ \boldsymbol{\omega}_{(n+1)} = (\boldsymbol{\omega}_{n+1}, \boldsymbol{\omega}_{n+2}, \cdots) \in \Omega_{(n+1)},$$

$$\mathcal{P}_{0,n} = \left\{ P_{|\mathcal{G}_n} : P \in \mathcal{P} \right\},$$

$$\mathcal{G}_{(n+1)} = \left\{ A : A \subset \Omega_{(n+1)}, \ \Omega^{(n)} \times A \in \mathcal{G} \right\}.$$

The set $\mathcal{P}_{0,n}$ represents the set of probability measures on the measurable space $(\Omega^{(n)}, \mathcal{G}_n)$. The regular conditional probability of a measure Q on the measurable space (Ω, \mathcal{G}) given the sub- σ -algebra \mathcal{G}_n is denoted by $Q(\cdot|\mathcal{G}_n)$. A probability kernel $\psi: \Omega^{(n)} \times \mathcal{G}_{(n+1)} \to [0,1]$ is called a \mathcal{P} -kernel if, for every n and every $\omega^{(n)} \in \Omega^{(n)}$, there exists a probability measure $Q \in \mathcal{P}$ such that

$$\psi(\boldsymbol{\omega}^{(n)}, A) = Q(\Omega^{(n)} \times A | \mathcal{G}_n)(\boldsymbol{\omega}^{(n)})$$
(3.5)

for all $A \in \mathcal{G}_{(n+1)}$.

A probability measure P on (Ω, \mathcal{G}) is said to be *rectangular* with respect to the filtration \mathcal{G}_n if, for every P, every P, every P, and every P-kernel P as above, the measure P defined by

$$P(A) = \int_{\Omega^{(n)}} \int_{\Omega_{(n+1)}} I_A\left(\boldsymbol{\omega}^{(n)}, \boldsymbol{\omega}_{(n+1)}\right) \psi\left(\boldsymbol{\omega}^{(n)}, d\boldsymbol{\omega}_{(n+1)}\right) P_n\left(d\boldsymbol{\omega}^{(n)}\right), \quad \forall A \in \mathcal{G},$$
 (3.6)

belongs to \mathcal{P} .

If the probability set \mathcal{P} consists of a single probability distribution P, then the principle of rectangularity is a straightforward consequence of the Bayesian updating rule. This is because P can be decomposed into its marginal and conditional distributions, which can then be recombined to obtain P again. However, when dealing with more complex sets \mathcal{P} , rectangularity is a non-trivial requirement. It demands that \mathcal{P} is closed under the operation of pasting together conditionals and marginals that may have been induced by different measures within \mathcal{P} . The significance of rectangularity is highlighted in the reference [6].

Specialize the above framework by assuming that there exists a measurable space $(\overline{\Omega}, \overline{\mathcal{F}})$ such that, for all $1 \leq i \leq n$, $(\Omega_i, \mathcal{F}_i) = (\overline{\Omega}, \overline{\mathcal{F}})$ and $\mathcal{G}_n = \prod_{i=1}^n \mathcal{F}_i$. That is, all positions in a sequence have a common set of possible outcomes $\overline{\Omega}$ and an associated common σ -algebra $\overline{\mathcal{F}}$. As an example, for a language modeling task, $\overline{\Omega}$ could be the set V of all possible tokens in a given language, i.e. $\overline{\Omega} = V$, and $\overline{\mathcal{F}}$ could be the collection of all subsets of $\overline{\Omega}$ that can be generated by using the given grammar rules and vocabulary V of the language.

We can define the *one-step-ahead conditionals* as follows: for each $P \in \mathcal{P}$ and each n, we let $P_{n,n+1}(\omega^{(n)})$ be the restriction to \mathcal{G}_{n+1} of $P(\cdot|\mathcal{G}_n)(\omega^{(n)})$. Assuming that there exists a subset \mathcal{L} of $\Delta(\overline{\Omega},\overline{\mathcal{F}})$, all of whose measures are equivalent, we can specialize the above framework by defining the indistinguishably and independently distributed (IID) model through the set \mathcal{P}^{IID} . Specifically,

$$\mathcal{P}^{\text{IID}} = \left\{ P \in \Delta(\Omega, \mathcal{G}) : P_{n,n+1}(\omega^{(n)}) \in \mathcal{L}, \ \forall n, \ \forall \omega^{(n)} \in \Omega^{(n)} \right\}. \tag{3.7}$$

The set \mathcal{P}^{IID} is comprised of measures for which the one-step-ahead conditionals, at every history, are in \mathcal{L} . In other words, \mathcal{L} represents the plausible probability laws for every position in a sequence, reflecting our partial ignorance about each position individually, regardless of its history.

The set \mathcal{P}^{IID} only requires that each one-step-ahead conditional belongs to the set of plausible probability laws \mathcal{L} , without imposing any further constraints on the measures of the entire sequence. It models partial ignorance about each position separately without considering the history.

However, it is possible to have a more complex situation where the one-step-ahead conditionals at each time point/position have different sets of probabilities, denoted by \mathcal{L}_n . In this case, the set

 \mathcal{P}^{LM} is defined similarly to \mathcal{P}^{IID} , but with the requirement that each one-step-ahead conditional belongs to the corresponding set of plausible probability laws \mathcal{L}_n . That is,

$$\mathcal{P}^{LM} = \left\{ P \in \Delta(\Omega, \mathcal{G}) : P_{n,n+1}(\omega^{(n)}) \in \mathcal{L}_n, \ \forall n, \ \forall \omega^{(n)} \in \Omega^{(n)} \right\}. \tag{3.8}$$

Our GAMMT model captures the probabilistic structure of uncertainty by utilizing the set \mathcal{P}^{LM} . Specifically, the GAMMT model allows for the possibility of different sets of plausible probability laws for each one-step-ahead conditional at every time point, as denoted by \mathcal{L}_n . By incorporating this more flexible set of probability laws, the GAMMT model is able to capture a wider range of uncertainty and account for potential dependencies between different positions in a sequence.

4 Architecture

In this section, we propose a class of deep neural networks that capture ambiguity in the data-generating process. Specifically, we introduce the GAMMT model, which utilizes multiple self-attention Transformer decoders commonly used in natural language processing tasks (e.g., [3, 10, 22, 23, 24, 26, 32]).

The architecture of the GAMMT model, shown in Figure 1, consists of M parallel Transformer decoders, each with N identical layers. Each layer includes a masked multi-head self-attention sublayer with h parallel attention heads, a feed-forward network, and layer normalization. The output of each decoder is transformed into a set of next-token probabilities, denoted by \mathcal{L}_n for the n-th position in a sequence, using a linear transformation and softmax function. The M Transformer decoders are connected by a selection mechanism S_n that captures the ambiguity in the data-generating process.

By utilizing the set $\mathcal{P}^{\mathrm{LM}}$ of plausible probability laws, the GAMMT model captures a wide range of uncertainty and potential dependencies between different positions in a sequence. The precise probabilistic structure of the GAMMT model is described in Section 3.2. The model learns the conditional probabilities at each time point n, which are then combined using the selection mechanism S_n to capture the ambiguity in the data-generating process.

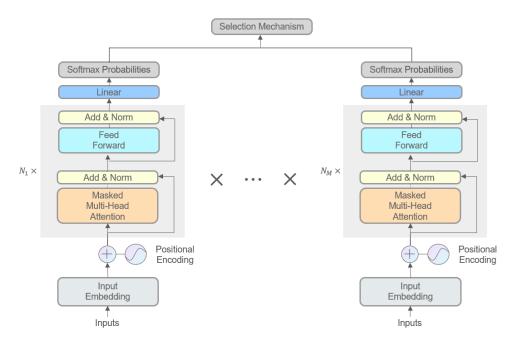


Figure 1: Model architecture - The Transformer decoders to approximate sets of probabilities

5 Algorithms

This section provides an overview of the algorithms for the proposed GAMMT model's architecture, model training, and inference. These algorithms are presented to help readers understand the workings of the GAMMT model. The pseudocode for the algorithms described in this section can be found in Appendix A.

Algorithm 1 outlines the forward pass of the model, which generates a set of probability matrices for predicting the next token, given an input sequence of tokens x with length ℓ_x . The vocabulary of tokens is denoted by V, with N_V being the size of the vocabulary. The operations of embedding and Transformer decoding are similar to those described in Algorithms 1, 2, and 10 of Phuong and Hutter's work [22]. Algorithm 2 presents the training process for the GAMMT model. The M parallel Transformer decoders are trained simultaneously using the given selection mechanism for N_{epochs} epochs to obtain the trained parameters. Algorithm 3 outlines the process of generating model predictions for a prompt using the GAMMT model, with the specific sampling method dependent on the selection mechanism used during training.

For brevity, we omit the details of the initial token and positional embedding, as well as the specifics of the Transformer decoders. Interested readers can refer to Phuong and Hutter's work [22] for a comprehensive overview of Transformer models and related approaches.

6 Conclusion

In conclusion, this paper presents a novel approach for modeling data generation processes under ambiguity, called GAMMT. Our proposed approach acknowledges that the process may be ambiguous and determined by multiple probabilities, which is different from traditional machine learning methods that assume data is sampled from a deterministic probability. GAMMT utilizes multiple parallel transformers to model a set of probabilities in sequential data, providing multiple representations of the underlying input tokens and sequence.

By leveraging the effectiveness of transformers, our GAMMT model has two key features that distinguish it from existing approaches. Firstly, it is generative and can effectively model ambiguity in data generation processes. Secondly, the last hidden layer of each parallel transformer provides multiple embeddings of each token and the input sequence, enabling diverse representations of the underlying input. We anticipate that our models will improve the quality and diversity of generated text, making it more engaging and human-like. Although we have not yet provided experimental validation, we believe that our models have great potential in this regard.

Our approach to modeling ambiguity in sequential data is novel and provides a more interpretable and flexible method than existing techniques. The unique features of our model, such as the use of multiple transformers and its generative nature, make it a promising direction for future research in machine learning. We hope that this framework will stimulate further research in the fields of machine learning and natural language processing, and that our work will contribute to the development of more effective and flexible models for handling ambiguity in data generation processes.

References

- [1] Amatriain X. Transformer models: an introduction and catalog. arXiv preprint arXiv:2302.07730, 2023.
- [2] Bach D R, Hulme O, Penny W D, et al. The known unknowns: neural representation of second-order uncertainty, and ambiguity. Journal of Neuroscience, 2011, 31(13): 4811-4820.
- [3] Brown T, Mann B, Ryder N, et al. Language models are few-shot learners. NeurIPS, 2020, 33: 1877-1901.
- [4] Buisson M, Alonso-Jiménez P, Bogdanov D. Ambiguity Modelling with Label Distribution Learning for Music Classification. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2022: 611-615.
- [5] Camerer C, Weber M. Recent developments in modeling preferences: Uncertainty and ambiguity. Journal of risk and uncertainty, 1992, 5(4): 325-370.
- [6] Chen Z, Epstein L G. A central limit theorem for sets of probability measures. Stochastic Processes and their Applications, 2022, 152: 424-451.

- [7] Chen M, Tworek J, Jun H, et al. Evaluating large language models trained on code. arXiv preprint arXiv:2107.03374, 2021.
- [8] Chowdhery A, Narang S, Devlin J, et al. PaLM: Scaling language modeling with pathways. arXiv preprint arXiv:2204.02311, 2022.
- [9] Chumbley J R, Flandin G, Bach D R, et al. Learning and generalization under ambiguity: an fMRI study. PLoS Computational Biology, 2012, 8(1): e1002346.
- [10] Devlin J, Chang M W, Lee K, et al. BERT: Pre-training of deep bidirectional transformers for language understanding. arXiv:1810.04805, 2018. NAACL, 2019.
- [11] Dosovitskiy A, Beyer L, Kolesnikov A, et al. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. arXiv:2010.11929, 2020.
- [12] Ek C H, Rihan J, Torr P H S, et al. Ambiguity modeling in latent spaces. International workshop on Machine Learning for Multimodal Interaction. Springer, Berlin, Heidelberg, 2008: 62-73.
- [13] Ellsberg, D. Risk, Ambiguity, and the Savage Axioms, Quarterly Journal of Economics, 1961, 75, 643-669.
- [14] Fedus W, Zoph B, Shazeer N. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. The Journal of Machine Learning Research, 2022, 23(1): 5232-5270.
- [15] Hoffmann J, Borgeaud S, Mensch A, et al. Training compute-optimal large language models. arXiv preprint arXiv:2203.15556, 2022.
- [16] Ilut C L, Schneider M. Modeling uncertainty as ambiguity: A review. NBER Working Paper, 2022.
- [17] Jumper J, Evans R, Pritzel A, et al. Highly accurate protein structure prediction with AlphaFold. Nature 596, 583-589, 2021.
- [18] Levy I, Snell J, Nelson A J, et al. Neural representation of subjective value under risk and ambiguity. Journal of neurophysiology, 2010, 103(2): 1036-1047.
- [19] Ouyang L, Wu J, Jiang X, et al. Training language models to follow instructions with human feedback[J]. Advances in Neural Information Processing Systems, 2022, 35: 27730-27744.
- [20] OpenAI (2022), ChatGPT: https://openai.com/blog/chatgpt
- [21] Patel R, Nenkova A. Modeling ambiguity in text: A corpus of legal literature. 2019.
- [22] Phuong M, Hutter M. Formal algorithms for transformers. arXiv:2207.09238, 2022.

- [23] Radford A, Narasimhan K, Salimans T, et al. Improving language understanding by generative pre-training. 2018.
- [24] Radford A, Wu J, Child R, et al. Language models are unsupervised multitask learners. OpenAI, 2019.
- [25] Rae J W, Borgeaud S, Cai T, et al. Scaling language models: Methods, analysis & insights from training gopher. arXiv preprint arXiv:2112.11446, 2021.
- [26] Raffel C, Shazeer N, Roberts A, et al. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. arXiv:1910.10683, 2019.
- [27] Scao T L, Fan A, Akiki C, et al. BLOOM: A 176B-parameter open-access multilingual language model. arXiv preprint arXiv:2211.05100, 2022.
- [28] Taya F. Seeking ambiguity: a review on neuroimaing studies on decision making under ambiguity. Recherches Économiques de Louvain/Louvain Economic Review, 2012, 78(03-04): 85-100.
- [29] Taylor R, Kardas M, Cucurull G, et al. Galactica: A large language model for science. arXiv preprint arXiv:2211.09085, 2022.
- [30] Thoppilan R, De Freitas D, Hall J, et al. LaMDA: Language models for dialog applications. arXiv preprint arXiv:2201.08239, 2022.
- [31] Touvron H, Lavril T, Izacard G, et al. LLaMA: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971, 2023.
- [32] Vaswani A, Shazeer N, Parmar N, et al. Attention is all you need. NeurIPS, 2017, 30.
- [33] Wei J, Bosma M, Zhao V Y, et al. Finetuned language models are zero-shot learners. arXiv preprint arXiv:2109.01652, 2021.
- [34] Yang Z, Dai Z, Yang Y, et al. Xlnet: Generalized autoregressive pretraining for language understanding. Advances in neural information processing systems, 2019, 32.
- [35] Yang G, Zhang J, Zhang Y, et al. Probabilistic modeling of semantic ambiguity for scene graph generation. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2021: 12527-12536.
- [36] Zeng A, Liu X, Du Z, et al. GLM-130B: An open bilingual pre-trained model. arXiv preprint arXiv:2210.02414, 2022.
- [37] Zhang S, Roller S, Goyal N, et al. OPT: Open pre-trained transformer language models. arXiv preprint arXiv:2205.01068, 2022.

Appendix A Algorithms

This appendix outlines the algorithms used for the proposed method, including model architecture, training, and inference. The GAMMT model architecture, which details the forward pass of the model, is presented in Algorithm 1. The model takes an input sequence of tokens x with a length of ℓ_x , and outputs a set of probability matrices for predicting the next token. The vocabulary of tokens is denoted by V, with N_V representing the vocabulary size. The embedding and Transformer decoding operations are similar to those described in Phuong and Hutter's work [22], as outlined in Algorithms 1, 2, and 10 ([22]).

```
Algorithm 1: \mathcal{P} \leftarrow \text{DTransformers}(x|\theta)
  /* Model architecture - the Transformer decoders, forward pass
  /* Input: a sequence of token IDs with length \ell_x
                                                                                                           */
  /* Output: probabilities of next token
  Input: x \in V^*
  Output: P_{\theta_m} \in [0,1]^{N_V \times \ell_x}, m \in [M] := \{1,2,\cdots,M\} and the Selection S
  Hyperparameters: The number of Transformer decoders M, the Transformer decoders
                          Layers N_i, heads H_i, input embedding dimension d_e, hidden layer
                          dimension d_{mlp}, maximum length of input \ell_{max}, selection S
  Parameters: \theta = (\theta_m)_{m \in [M]} includes all token and positional embedding matrices, all
                  Transformers' parameters
1 for m = 1, 2, \dots, M do
      X_m \leftarrow \text{TransformerDecoder}(m, \text{Embedding}(x))
    P_{\theta_m} = \operatorname{softmax}(W_o^{(m)} X_m)
4 end
s return (P_{\theta_m})_{m\in[M]}, S=S(P_{\theta_1},\cdots,P_{\theta_M})
```

The training process for the GAMMT model is presented in Algorithm 2. The model employs M parallel Transformer decoders, trained simultaneously with a selection mechanism. The trained parameters are obtained after N_{epochs} epochs of training.

Algorithm 3 outlines the process of generating model predictions from a given prompt using the GAMMT model. The output is a continuation of the prompt generated by the model, with the sampling method used during inference depending on the selection mechanism employed during training.

```
Algorithm 3: y \leftarrow Inference(x, \hat{\theta})
   /* Model inference - generate a sequence based on the trained model and a prompt
         */
    /* Input: the trained parameters and a prompt
                                                                                                                                 */
    /* Output: a continuation of the prompt sampled from the trained model
                                                                                                                                 */
   Input: \hat{\theta} = (\hat{\theta}_m)_{m \in [M]}, x \in V^*
    Output: y \in V^*
   Hyperparameters: temprature \tau > 0
 1 \ell_x \leftarrow \text{length}(x)
 y \leftarrow \emptyset
 3 while y \neq eos\_token do
         (P_{\theta_m})_{m \in [M]}, -\leftarrow \text{DTransformers}(x|\hat{\theta})
         if S \sim \mathcal{R}([M]) then
 5
              u \leftarrow \text{sample a random variable } \mathcal{R} \text{ on } [M]
              p \leftarrow P_{\hat{\theta}_u}[:, \ell_x]
 7
              sample a token y from the probability q \propto p^{1/\tau}
              x \leftarrow [x, y]
 9
            \ell_x \leftarrow \ell_x + 1
10
         end
11
         else if S = \max then
12
              for m = 1, 2, \dots, M do
13
                   p_m \leftarrow P_{\hat{\theta}_m}[:, \ell_x]
14
                   sample a token y_m from the probability q_m \propto p_m^{1/\tau}
15
              end
16
              y \leftarrow y_{m^*} with m^* = \operatorname{argmax} \{q_m(y_m), m \in [M]\}
17
              x \leftarrow [x, y]
18
              \ell_x \leftarrow \ell_x + 1
19
         end
20
21 end
```

22 return y = x