# GAMMT: Generative Ambiguity Modeling Using Multiple Transformers

Xingcheng Xu

November 2022
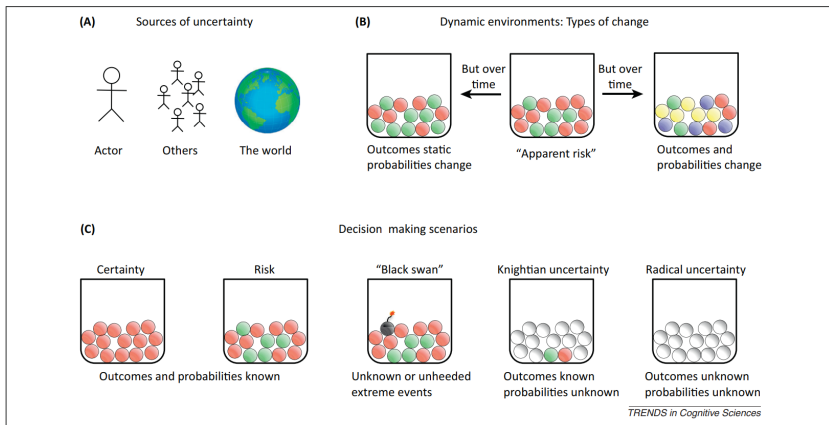
# Outline

# Introduction

# Questions

1. What is uncertainty(risk/ambiguity/...)? How much important is it?
2. Where are the sources of ambiguity from? How to understand and model ambiguity?
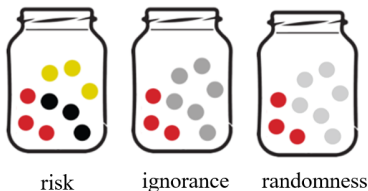3. What benefits can we get from recognizing the presence of ambiguity with respect to the pure risk situation?

Source: Meder et al. (2013), Trends in Cognitive Sciences.

# Examples of Uncertainty

- Urn 1: 30 red, 30 black and 30 yellow balls (*risk/deterministic probability*)
- Urn 2: 30 red balls and 60 black or yellow balls (*subjective ignorance*)
- Urn 3: 90 balls, 30 in red but the number of black or yellow balls are uniformly distributed in this magic urn (*objective randomness*)



risk    ignorance    randomness

# Risk and Ambiguity

- Risk refers to situations where the perceived likelihoods of events of interest can be represented by a probability measure.
- Ambiguity refers to situations where the information available to the decision-maker is too imprecise to be summarized by a probability measure (*subjective ignorance*) or results from genuine randomness in the external environment (*objective randomness*).

- **Ellsberg Paradox:** an urn containing 90 balls, identical except for color. You know that exactly 30 of the balls are red. Each of the remaining 60 balls is either black or yellow, but you do not know the relative numbers of black and yellow balls.[1]

|       | 30 | 60 | |
|       | red | black | yellow |
| ----- | ----- | ----- | ----- |
| $f_1$ | \$100 | \$0   | \$0   |
| $f_2$ | \$0   | \$100 | \$0   |
| $f_3$ | \$100 | \$0   | \$100 |
| $f_4$ | \$0   | \$100 | \$100 |

- **Preference:** empirically $f_1 \succ f_2$ but $f_4 \succ f_3$, which contradicts with the risk-based models.

---

[1] Ellsberg (1961), Machina and Schmeidler (1992)

# Ambiguity and Decision-Making

- Ellsberg Paradox demonstrated that the choices under ambiguity are distinct from risk and such a distinction is *behaviorally meaningful*, and suggests that ambiguity is at least as prominent as risk in making investment decisions.

- An agent using the wrong probability measure may plausibly be aware of this possibility and thus be led to seek robust decisions. Such self-awareness and a desire for robust decisions lead naturally to consideration of sets of priors.

- Different agent groups use different probability measures to generate contents or make decisions, leading naturally to consideration of sets of probabilities.

# Language, Image, ...

Generated by people (decision) or environment (semantic ambiguity, perceptual vagueness)

- incomplete knowledge, limited information (feeling, perception, experiences)
- interaction with heterogeneous agents
- genuine randomness in the external environment

# Literature

- Nature Language Processing:
  - ▶ Franz. (1996). Automatic ambiguity resolution in natural language processing: an empirical approach. Berlin, Heidelberg: Springer.
  - ▶ Anjali & Babu (2014). Ambiguities in natural language processing. International Journal of Innovative Research in Computer and Communication Engineering.
  - ▶ Patel & Nenkova (2019). Modeling ambiguity in text: A corpus of legal literature.
  - ▶ Jackson (2020). Understanding understanding and ambiguity in natural language. Procedia Computer Science.
  - ▶ Yadav, Patel & Shah (2021). A comprehensive review on resolving ambiguities in natural language processing. AI Open.

# Literature

- Computer Vision:
  - ▶ Casadei & Mitter (1999). Beyond the uniqueness assumption: Ambiguity representation and redundancy elimination in the computation of a covering sample of salient contour cycles. Computer Vision and Image Understanding.
  - ▶ Rust & Stocker (2010). Ambiguity and invariance: two fundamental challenges for visual processing. Current opinion in neurobiology.
  - ▶ Rupprecht et al. (2017). Learning in an uncertain world: Representing ambiguity through multiple hypotheses. ICCV.
  - ▶ Manhardt et al. (2019). Explaining the ambiguity of object detection and 6d pose from visual data. ICCV.
  - ▶ Yang et al. (2021). Probabilistic modeling of semantic ambiguity for scene graph generation. CVPR.

# Literature

- Science:
  - ▶ Borodovsky & McIninch (1993). Recognition of genes in DNA sequence with ambiguities. Biosystems.
  - ▶ Postic et al. (2017). An ambiguity principle for assigning protein structural domains. Science advances.
  - ▶ Barbieri (2019). Evolution of the genetic code: The ambiguity-reduction theory. Biosystems.
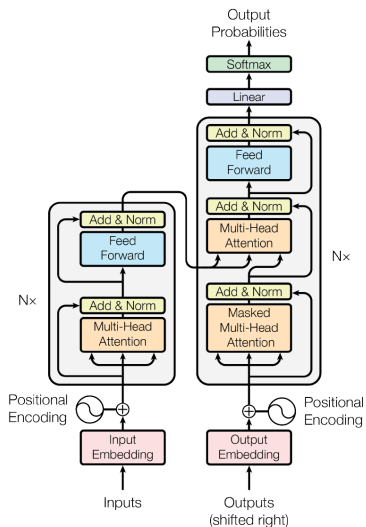  - ▶ Oliva et al. (2019). Accounting for ambiguity in ancestral sequence reconstruction. Bioinformatics.
- NeuroScience:
  - ▶ Levy et al. (2010). Neural representation of subjective value under risk and ambiguity. Journal of neurophysiology.
  - ▶ Bach et al. (2011). The known unknowns: neural representation of second-order uncertainty, and ambiguity. Journal of Neuroscience.
  - ▶ Chumbley et al. (2012). Learning and generalization under ambiguity: an fMRI study. PLoS Computational Biology.

# The Transformer

# Transformers

- The landmark paper:
  - The Transformer (Vaswani et al., 2017)
- Nature Language Processing:
  - BERT (Devlin et al., 2018),
  - OpenAI GPT (Radford et al., 2018), GPT-2 (Radford et al., 2019), GPT-3 (Brown et al., 2020),
  - Google T5 (Raffel et al., 2019), ...
- Computer Vision:
  - Vision Transformer (Dosovitskiy et al., 2020), ...
- Science: Structural Biology
  - AlphaFold 2 (Jumper et al., 2021), ...
- AIGC:
  - CLIP (Radford et al., 2021),
  - DALL-E 2 (Ramesh et al., 2022), Stable Diffusion (Rombach et al., 2022), eDiff-I (Balaji et al., 2022), ...

Source: Vaswani et al. (2017)

# The Transformer: a Probability Machine



Source: Alammar (2018), The Illustrated Transformer

$$P(s_t | s_1, s_2, \cdots, s_{t-1})$$

# Generative Ambiguity Modeling

## Conventional Wisdom in AI

Standard approaches in machine learning for sequential modeling such as language models are to factorize the joint probabilities $P(x)$ of a sequence $x = (s_1, s_2, \cdots, s_n) \in V^*$ in a space of sequences $V^* = \cup_{\ell=0}^{\infty} V^\ell$ ($V$ is a vocabulary of tokens) as the product of conditional probabilities (e.g. GPT, GPT-2, GPT-3)

$$P(x) = \prod_{t=1}^{n} P(s_t | s_1, s_2, \cdots, s_{t-1}) \tag{1}$$

or to factorize the conditional probabilities $P(x|z)$ of a pair of sequences $x = (s_1, s_2, \cdots, s_n) \in V^*$ and $z \in V^*$ as the product of following conditionals (e.g. the original Transformer)

$$P(x|z) = \prod_{t=1}^{n} P(s_t | s_1, s_2, \cdots, s_{t-1}, z). \tag{2}$$

# Sequence With Ambiguity

- $(\Omega, \{\mathcal{G}_i\}_{i=1}^\infty, \mathcal{G}) = (\prod_{i=1}^\infty \Omega_i, \{\mathcal{G}_i\}_{i=1}^\infty, \sigma(\cup_{i=1}^\infty \mathcal{G}_i))$ a filtered space modeling a sequence of experiments/games/events/objects.

- The set of possible outcomes for the $i^{th}$ experiment is $\Omega_i$. $\mathcal{G}_i$ is a $\sigma$-algebra on $\prod_{j=1}^i \Omega_j$ representing the information regarding experiments $1, 2, \cdots, i$.

- The ex ante probabilities of experiments are not known precisely and are represented by a set $\mathscr{P}$ of probability measures on $(\Omega, \mathcal{G})$, and assume that all measures in $\mathscr{P}$ are equivalent on each $\mathcal{G}_n$.

# Probability Structure: Rectangularity

- $\mathscr{P}_{0,n} := \{P_{|\mathcal{G}_n} : P \in \mathscr{P}\}$, $\Omega = \Omega^{(n)} \times \Omega_{(n+1)} = \prod_{i=1}^{n} \Omega_i \times \prod_{i=n+1}^{\infty} \Omega_i$

- Probability kernel: a functional $\lambda : \Omega^{(n)} \times \mathcal{G}_{(n+1)} \to [0,1]$ satisfying $\mathcal{G}_n$-measurable and being a probability measure on $(\Omega_{(n+1)}, \mathcal{G}_{(n+1)})$.

- $\mathscr{P}$-kernel: If $\forall \omega^{(n)} \in \Omega^{(n)}$, $\exists Q \in \mathscr{P}$ satisfying

$$\lambda(\omega^{(n)}, A) = Q(\Omega^{(n)} \times A | \mathcal{G}_n)(\omega^{(n)}), \ \forall A \in \mathcal{G}_{(n+1)}.$$

- $\mathscr{P}$ is *rectangular* if $\forall n \in \mathbb{N}$, $\forall p_n \in \mathscr{P}_{0,n}$ and for every $\mathscr{P}$-kernel $\lambda$, if $P$ is defined as

$$P(A) := \int I_A(\omega) \lambda(\omega^{(n)}, d\omega_{(n+1)}) p_n(d\omega^{(n)}), \ \forall A \in \mathcal{G}, \qquad (3)$$

then $P \in \mathscr{P}$.

- $\mathscr{P}$ is closed w.r.t. pasting of alien marginals and conditionals, endowing $\mathscr{P}$ with a recursive structure.

# Example of Rectangularity: IID Model

- Experiments have a common set of possible outcomes $\overline{\Omega}$ and a common $\sigma$-algebra $\overline{\mathcal{F}}$.
- Fix a subset $\mathcal{L} \subset \mathscr{M}(\overline{\Omega}, \overline{\mathcal{F}})$ & all measures in $\mathcal{L}$ are equivalent.
- Let $P_{n,n+1}(\omega^{(n)})$ denote the restriction to $\mathcal{G}_{n+1}$ of $P(\cdot|\mathcal{G}_n)(\omega^{(n)})$.
- IID (**I**ndistinguishably and **I**ndependently **D**istributed) model:

$$\mathscr{P}^{IID} = \left\{ P \in \mathscr{M}(\Omega, \mathcal{G}) : \ P_{n,n+1}(\omega^{(n)}) \in \mathcal{L}, \ \forall n, \omega^{(n)} \in \Omega^{(n)} \right\}$$

- The set consists of all measures whose *one-step-ahead conditionals*, at every history, lie in $\mathcal{L}$ modeling partial ignorance about each experiment separately.
- If $\mathcal{L} = \{P\}$, $\mathscr{P}^{IID}$ degenerates to a probability (the conventional probability model, random walk).

# Example of IID Model

- Each experiment can produce one of the three outcomes (3 words in a vocabulary): $\overline{\Omega} = \{w_1, w_2, w_3\}$ and $\overline{\mathcal{F}} = 2^{\overline{\Omega}}$.

- Probabilities are not known exactly but it is known that, for each experiment, the outcomes are given by ($0 < q < p,\ p + q \le 1$)

$$\mathcal{L} = \{P_1 = (p, q, 1 - p - q), P_2 = (q, p, 1 - p - q)\}.$$

- The ignorance about the relation between experiments is subject to the IID model $\mathscr{P}^{IID}$.

# Deep Ambiguity Sequence Model: GAMMT

# Deep Ambiguity Sequence Model: GAMMT
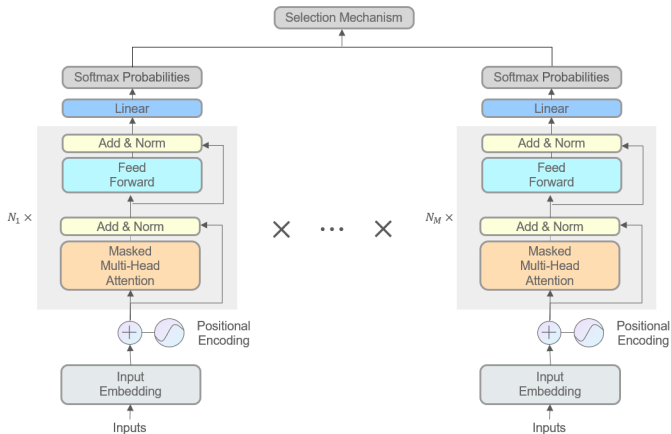
The GAMMT: An Ambiguity Machine



Figure: Model architecture - The Transformer decoders to approximate a set of probabilities

# Deep Ambiguity Sequence Model: GAMMT

There is a set of *one-step-ahead conditional probabilities* produced by the GAMMT model at every time step/every position of a sequence. The size of the set characterizes the level of ambiguity of data.

GAMMT: $\mathcal{P}(x) = \{P(x)|\ P \in \mathcal{M}\}$, $\mathcal{M} = \{P_{\theta_1}, P_{\theta_2}, \cdots, P_{\theta_M}\}$



$P_{\theta_m}(s_t|s_1, s_2, \cdots, s_{t-1}), m \in [M]$

# Deep Ambiguity Sequence Model: GAMMT

The model described above can be generalised to the encoder-decoder architectures with Transformers to approximate sets of conditional probabilities

$$\mathcal{P}(x|z) = \{P(x|z)| \ P \in \mathcal{M}\}$$

where $z$ and $x$ are a pair of input (source) or output (target) token sequences as the original Transformer/Google T5. (Could share the same encoder, or every parallel transformer has its only encoder.)

# Algorithm: Architecture

---

**Algorithm 1:** $\mathcal{P} \leftarrow \text{DTransformers}(x|\theta)$

---

```
/* Model architecture - the Transformer decoders, forward pass   */
/* Input: a sequence of token IDs with length ℓₓ                 */
/* 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 $N_j$, heads $H_j$, 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** $m = 1, 2, \cdots, M$ **do**

2      $X_m \leftarrow \text{TransformerDecoder}(m, \text{Embedding}(x))$

3      $P_{\theta_m} = \text{softmax}(W_o^{(m)} X_m)$

4 **end**

5 **return** $(P_{\theta_m})_{m \in [M]}$, $S = S(P_{\theta_1}, \cdots, P_{\theta_M})$

---

# Algorithm: Training

---

**Algorithm 2:** $\hat{\theta} \leftarrow \text{Training}(x_{1:N_{data}}, \theta)$

/* Model training - prediction of the next token                                    */
/* Input:  a dataset of sequences and initial parameters                            */
/* Output:  the trained parameters                                                  */

**Input:** $\{x_n\}_{n=1}^{N_{data}}$, $\theta = (\theta_m)_{m \in [M]}$

**Output:** $\hat{\theta} = (\hat{\theta}_m)_{m \in [M]}$

**Hyperparameters:** Epochs $N_{epochs}$, learning rate $\eta > 0$ and the selection mechanism $S$

1   **for** $i = 1, 2, \cdots, N_{epochs}$ **do**

2     **for** $n = 1, 2, \cdots, N_{data}$ **do**

3       $\ell_x \leftarrow \text{length}(x_n)$

4       $(P_{\theta_m})_{m \in [M]}, \_ \leftarrow \text{DTransformers}(x_n | \theta)$

5       $\text{loss}(\theta) = -\sum_{t=1}^{\ell_x - 1} \log S\left(\{P_{\theta_m}[x_n[t+1], t]\}_{m=1}^{M}\right)$

6       $\theta \leftarrow \theta - \eta \cdot \nabla \text{loss}(\theta)$

7     **end**

8   **end**

9   **return** $\hat{\theta} = \theta$

---

# Algorithm: Inference

---

**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  $\ell_x \leftarrow$ length$(x)$
2  $y \leftarrow \emptyset$
3  **while** $y \neq eos\_token$ **do**
4       $(P_{\theta_m})_{m \in [M]}, \_ \leftarrow$ DTransformers$(x|\hat{\theta})$
5       **if** $S \sim \mathcal{R}([M])$ **then**
6           $u \leftarrow$ sample a random variable $\mathcal{R}$ on $[M]$
7           $p \leftarrow P_{\hat{\theta}_u}[:, \ell_x]$
8           sample a token $y$ from the probability $q \propto p^{1/\tau}$
9           $x \leftarrow [x, y]$
10          $\ell_x \leftarrow \ell_x + 1$
11      **end**
12      **else if** $S =$ max **then**
13          **for** $m = 1, 2, \cdots, M$ **do**
14              $p_m \leftarrow P_{\hat{\theta}_m}[:, \ell_x]$
15              sample a token $y_m$ from the probability $q_m \propto p_m^{1/\tau}$
16          **end**
17          $y \leftarrow y_{m^*}$ with $m^* = \text{argmax}\{q_m(y_m), m \in [M]\}$
18          $x \leftarrow [x, y]$
19          $\ell_x \leftarrow \ell_x + 1$
20      **end**
21 **end**
22 **return** $y = x$

# Validation on NLP Tasks

- Dataset
  - ▶ CommonCrawl
  - ▶ WebText
  - ▶ Books1, Books2
  - ▶ Wikipedia
- Tasks
  - ▶ Machine translation
  - ▶ Semantic similarity
  - ▶ Text classification
  - ▶ Reading comprehension, summarization
  - ▶ Natural language inference, textual entailment
  - ▶ Commonsense reasoning, Question answering
  - ▶ Cloze tasks, Sentence/paragraph completion tasks
  - ▶ Article generation

Ref. GPT, GPT-2, GPT-3, T5, ...

# More

- Computer Vision: image as patch token sequence
- Audio Generation
- Protein Structure Prediction
- DNA/RNA Structure Prediction
- Reinforcement Learning
- ...

# Concluding Remarks

- long range dependencies (attention)
- zero-shot learning (scale)
- admitting ambiguity (structure)
- high diversity of generation (structure)
- multiple representations (structure)
- parallel computing, high computational efficiency (structure)
- ...

Thanks!