# The Free Transformer: 论文完整翻译与深度解读报告

## 摘要

### 原文翻译

We propose an extension of the decoder Transformer that conditions its generative process on random latent variables which are learned without supervision thanks to a variational procedure. Experimental evaluations show that allowing such a conditioning translates into substantial improvements on downstream tasks.

我们提出了一种对解码器Transformer的扩展,它通过一个无需监督的变分过程学习到的随机潜变量来调节其生成过程。实验评估表明,允许这种调节能够在下游任务中带来显著的性能提升。

## 深度解读

这篇论文的摘要虽然只有两句话,但信息量巨大,为我们揭示了一种可能引领下一代大语言模型架构的新思路。让我们逐字逐句地拆解,让你能像一位AI研究者一样理解其深层含义。

首先,"解码器Transformer"(decoder Transformer)是你非常熟悉的模型,比如GPT系列。它们的核心工作方式是"自回归"(autoregressive),即根据已经生成的词语,来预测下一个最有可能的词语,就像玩成语接龙一样,一个接一个。

本文的核心创新在于引入了"随机潜变量"(random latent variables)来"调节其生成过程"。这里的"潜"(latent)字是关键,意味着"隐藏的、不可见的"。你可以把它想象成一位作家在动笔写小说之前,脑海中先构思好的一个"核心主题"或"写作大纲"。这个大纲(比如"写一个悲伤的爱情故事"或"写一篇充满悬疑的侦探小说")并不会直接出现在小说的文字里,但它却像一只无形的手,在幕后指导着作者的每一个遣词造句。对于AI模型来说,这个"潜变量"就是它在生成文本之前,在内部为自己设定的一个"隐藏计划"。这个计划是"随机"的,意味着模型可以自由地探索各种可能的计划。

那么,模型如何学会制定这些"计划"呢?答案是"无需监督的变分过程"(unsupervised variational procedure)。"无需监督"意味着我们不需要人工去告诉模型"这是一个悲伤的故事"或者"那是一篇科学论文"。我们只需要把海量的文本数据(比如整个互联网)喂给它,它就会通过一种名为"变分自编码器"(Variational Autoencoder, VAE)的数学框架,自己去学习和归纳文本中存在的各种潜在结构和主题。这就像一个孩子通过大量阅读,自己领悟到了不同文体和情感的写作技巧,而不是靠老师手把手地教。

最后,也是最重要的成果: "在下游任务中带来显著的性能提升"。"下游任务"指的是我们希望AI完成的具体工作,比如编写代码、解答数学题、进行常识推理等。这篇论文通过实验证明,让模型先在内部形成一个"隐藏计划",然后再进行生成,能够让它在处理这些复杂任务时表现得更好。这暗示着,从简单的"词语接龙"模式,进化到更具规划性的"先定大纲,再写文章"模式,可能是提升AI"思考"和"推理"能力的关键一步。这不仅仅是对现有模型的微小改进,而是一次对底层工作机制的深刻变革。

## 1. 引言

## 原文翻译

Since their invention, the Transformer (Vaswani et al., 2017), and more specifically the decoder-only Transformers used originally for the GPT series of models (Radford et al., 2018), have become the core components of Al systems.

自Transformer架构被发明以来(Vaswani et al., 2017),特别是最初用于GPT系列模型的纯解码器 Transformer(Radford et al., 2018),已成为人工智能系统的核心组成部分。

## 深度解读

这一段为整篇论文的研究背景进行了定位。它首先强调了Transformer架构,特别是"纯解码器"(decoder-only)版本,在当今AI领域的统治地位。当你使用任何一个主流的大语言模型时,其技术内核几乎都源于此。作者提及2017年的开创性论文《Attention Is All You Need》和2018年的GPT系列,是为了说明这个架构已经经受了近十年的考验,并取得了巨大的成功。这既是对前人工作的致敬,也为接下来的颠覆性提议埋下了伏笔:既然这个架构如此成功和稳固,为什么我们还需要改变它呢?

## 原文翻译

It is remarkable that, after almost a decade, and in spite of improvements on many aspects of this class of methods, the autoregressive modelling of Transformers remains essentially unchallenged. We propose in this paper to revisit this key design aspect by allowing richer and more natural density models to emerge:

值得注意的是,在近十年之后,尽管这类方法在许多方面都有所改进,但Transformer的自回归建模方式基本上没有受到挑战。在本文中,我们建议重新审视这一关键设计,通过允许更丰富、更自然的密度模型出现来做出改变:

## 深度解读

这里,作者指出了当前技术范式中的一个"盲点"。尽管模型规模、训练数据和算法细节都在不断优化,但最核心的"自回归建模"(autoregressive modelling)方式——即逐字生成——却像一个不容置疑的"公理"被保留了下来。作者认为,是时候挑战这个"公理"了。他们提出的目标是构建"更丰富、更自然的密度模型"。这里的"密度模型"是一个偏学术的说法,你可以简单理解为"对文本生成规律的建模方式"。作者认为,逐字生成这种方式虽然有效,但可能不够"丰富"和"自然",因为它无法很好地捕捉文本的全局结构。人类写作时,脑中有一个整体的构思,而不是单纯地思考下一个词。作者希望让模型也能模拟这种更"自然"的生成方式。

### 原文翻译

• We extend the auto-regressive model of the decoder Transformer by allowing the conditioning on latent variables, thanks to a formulation as a conditional Variational Autoencoder (§ 3.1).

- We propose an implementation that requires a very modest computational and memory usage overhead (§ 3.2).
- 我们通过将其构建为一个条件变分自编码器(conditional Variational Autoencoder)的范式,引入潜变量作为条件,从而扩展了纯解码器Transformer的自回归模型(§ 3.1)。
- 我们提出了一种实现方式,其所需的计算和内存开销都非常小(§ 3.2)。

## 深度解读

这两点是本文的核心贡献,也是实现上述目标的具体路径。第一点,再次强调了解决方案的核心:潜变量(latent variables)和条件变分自编码器(cVAE)。"条件"(conditioning)这个词非常重要,它意味着模型的生成过程不再仅仅依赖于前面的文字,还要依赖于一个额外的、隐藏的"指令"——也就是潜变量。cVAE就是实现这一机制的数学工具。第二点,则展示了这项研究的工程价值。"计算和内存开销都非常小"意味着这个新方法并非是一个停留在理论上的美好构想,而是可以被实际应用到现有的大模型中,且不会带来高昂的成本。在动辄需要数千块GPU进行训练的今天,这是一个极具吸引力的优点。它意味着这项技术有潜力被广泛部署,而不是成为一个昂贵的"实验室玩具"。

## 原文翻译

The benefits of the proposed method are shown by training 1.5B and 8B models from scratch and assessing performance on multiple downstream benchmarks (§ 4).

我们通过从头开始训练15亿(1.5B)和80亿(8B)参数的模型,并在多个下游基准测试上评估其性能,来展示所提出方法的好处(§ 4)。

### 深度解读

这是对研究方法和成果的预告。作者明确表示,他们不是在小规模的"玩具"模型上进行实验,而是直接挑战了业界公认的"中量级"(1.5B)和"重量级"(8B)模型。这大大增强了实验结果的说服力。"从头开始训练"(training from scratch)也很关键,这排除了新旧模型混合可能带来的干扰,确保了实验的公平性。最后,通过在"多个下游基准测试"上进行评估,作者承诺将从不同维度全面地展示新架构的优势,而不仅仅是在某个特定任务上表现出色。这预示着接下来的实验部分将会有详实的数据来支撑他们的论点。

# 2. 动机

## 原文翻译

Decoder Transformers are auto-regressive discrete density approximators. They model a sequence of tokens  $S_1, ..., S_T$  by estimating the conditional distribution of each given those preceding it. Sampling is done by generating one token after another, each time computing the distribution of the next symbol given those generated so far.

解码器Transformer是自回归的离散密度近似器。它们通过估计每个词元(token)在给定其前面所有词元下的条件分布,来对一个词元序列  $S_1,...,S_T$  进行建模。采样是通过逐个生成词元来完成

的,每次都基于已生成的词元计算下一个符号的分布。

## 深度解读

这一段用非常精确的语言定义了标准解码器Transformer的工作原理。让我们把它翻译成更易懂的语言。把它想象成一个超级智能的输入法。当你输入"今天天气真"时,它会做什么?它会计算接下来成千上万个汉字出现的概率,可能会给"好"一个很高的概率,给"不错"一个次高的概率,而给"差"一个很低的概率。这就是"估计条件分布"。然后,它从这个概率分布中"采样"一个词,比如"好",并将其添加到句子末尾。接着,它基于"今天天气真好",继续预测下一个词。这个"一个接一个"的过程,就是"自回归"(auto-regressive)。"离散密度近似器"指的是它处理的是离散的单元(词元,tokens),并且试图学习这些词元组合在一起的概率规律。

### 原文翻译

The only density modelling and sampling that such models implement is that of the generated tokens. In particular, a decoder Transformer does not make additional latent decisions about the stream of symbols to generate. Its only decisions are the choices of the tokens themselves.

这类模型所实现的唯一密度建模和采样,仅仅是针对生成的词元本身。特别地,解码器Transformer 不会对要生成的符号流做出额外的潜在决策。它唯一的决策就是选择词元本身。

### 深度解读

这是对传统模型局限性的精准剖析。作者指出,标准模型的"决策空间"非常狭窄,它每一步唯一能做的决定就是"下一个词是什么?"。它缺乏一种更高层次的、全局性的决策能力。换句话说,它是一个战术家,精通于每一步的小规模战斗(选择下一个词),但它不是一个战略家,无法在战斗开始前就规划好整个战役的走向。这种模式缺乏"潜在决策"(latent decisions),也就是我们之前提到的"隐藏计划"。

### 原文翻译

Consider, for instance, that we train such a model to generate movie reviews and that we want to have two clearly separated categories of negative and positive reviews. Given a large enough model and the necessary amount of training data, there is no doubt that a decoder Transformer trained on a dataset of that form would work perfectly and would generate these two types of reviews. However, to do so, it would generate tokens one after another and decide, based on the words generated so far, whether the review it is currently generating is a positive or a negative one, and continue the process accordingly. In particular, the model would not make the explicit decision to generate a negative or a positive review. It would produce tokens, and this notion of a negative or positive review would be implicit in their posterior probabilities.

举个例子,假设我们训练这样一个模型来生成电影评论,并且我们希望有正面和负面两类清晰分开的评论。只要模型足够大,训练数据足够多,毫无疑问,一个在这种数据集上训练的解码器 Transformer能够完美地工作,并生成这两种类型的评论。然而,为了做到这一点,它会逐个生成词元,并根据已经生成的词语来判断它当前正在生成的评论是正面的还是负面的,然后相应地继续这 个过程。特别地,模型不会做出一个明确的决定去生成一篇负面或正面的评论。它只是产出词元,而这种正面或负面的概念是隐含在这些词元的后验概率中的。

## 深度解读

这个电影评论的例子非常精彩,它生动地揭示了标准模型工作方式的"笨拙"之处。 一个标准的大模型当然可以写出优秀的正面或负面影评。但它的"思考"过程是这样的:

- 1. 开始写第一个词,比如"这部"。
- 2. 接下来,它可能会生成"电影"。
- 3. 然后,基于"这部电影",它可能会生成"真是"。
- 4. 现在,基于"这部电影真是",模型内部的概率开始发生变化。如果它接下来生成了"太棒了",那么整个句子的"正面"倾向就大大增加,后续它就会更倾向于选择"精彩"、"推荐"等词。如果它生成了"太糟糕了",那么"负面"倾向就会增加,后续就会更倾向于选择"失望"、"浪费时间"等词。

看到了吗?模型是在"摸着石头过河"。它通过已经生成的词语,反过来推断自己"应该"在写哪种类型的评论。它没有一个"上帝视角"的预设。作者指出,这种"正面或负面"的概念,是"隐含在后验概率中的",而不是一个预先做出的"明确决定"。这种方式不仅效率低下,而且容易出错。如果模型在开头不小心用了一个模棱两可的词,就可能导致后续的生成内容精神分裂,前后矛盾。

相比之下,Free Transformer的工作方式更像人类:

- 1. 在动笔前,先在内部设定一个潜变量 Z= "写一篇负面评论"。
- 2. 然后,整个生成过程都以这个"指令"为条件。模型会直接从"负面词汇库"和"负面句式库"中进行选择,生成诸如"这部电影真是太糟糕了,完全是浪费时间"这样的评论。 这种方式更加直接、高效,且保证了全局的一致性。

## 原文翻译

Due to the chain rule, any density can be modelled as autoregressive. However, in particular when the "natural" structure involves conditioning on latent variables, the autoregressive model of the signal may be a great deal more complex than the full joint model including the latent.

根据链式法则,任何密度都可以被建模为自回归形式。然而,特别是当"自然的"结构涉及到对潜变量的条件依赖时,信号的自回归模型可能会比包含该潜变量的完整联合模型复杂得多。

### 深度解读

这一段从数学原理上解释了为什么传统模型会那么"笨拙"。"链式法则"是概率论中的一个基本定理,它允许我们将一个复杂的联合概率分解成一系列条件概率的乘积。这正是自回归模型能够"理论上"模拟任何文本分布的数学基础。但是,作者指出了一个关键问题:理论上可行,不代表实践中高效。当一个现象的"自然结构"本身就包含一个核心的隐藏变量时(比如影评的"正面/负面"属性),强行用纯自回归模型去模拟它,会导致模型本身变得"异常复杂"。这就好比,我们知道地球围绕太阳转(一个简单的、包含隐藏中心"太阳"的模型),但如果你坚持以地球为宇宙中心来描

述天体运动,你就需要建立一套极其复杂的"本轮-均轮"体系来解释行星的逆行。虽然也能解释,但模型变得复杂、丑陋且不自然。

## 原文翻译

We can consider a simple example illustrating that point. Let  $Z \sim B(0.5)$  be a latent "coin flip", and  $X_1, ..., X_T$  be equal to Z with independent flips of probability  $\epsilon$ . The  $X_{tS}$  are conditionally independent given Z, and we have  $P(X_{t+1} = 1 \mid Z = z) = \epsilon z + (1 - \epsilon)(1 - z)$  (1) however, expressed as an auto-regressive model without Z, we get:  $P(X_{t+1} = 1 \mid X_1 = z) = \epsilon z$ 

$$X_{1},...,X_{t}=X_{t}=\frac{\frac{\left(\frac{\epsilon}{1-\epsilon}\right)^{\sum_{j=1}^{t}x_{s}}(1-\epsilon)^{t+1}+\left(\frac{1-\epsilon}{\epsilon}\right)^{\sum_{s=1}^{t}x_{s}}\epsilon^{t+1}}{\left(\frac{\epsilon}{1-\epsilon}\right)^{\sum_{j=1}^{t}x_{s}}(1-\epsilon)^{t}+\left(\frac{1-\epsilon}{\epsilon}\right)^{\sum_{j=1}^{t}x_{s}}\epsilon^{t}}}$$
(2)

我们可以考虑一个简单的例子来说明这一点。让  $Z \sim B(0.5)$  是一个潜在的"抛硬币"结果,而  $X_1,...,X_7$  等于Z,但有  $\epsilon$  的概率会独立地翻转。  $X_t$  在给定Z的条件下是独立的,我们有  $P(X_{t+1}=1\mid Z=z)=\epsilon z+(1-\epsilon)(1-z)$  (1) 然而,如果表示成一个没有Z的自回归模型,我们得到:

$$P(X_{t+1} = 1 \mid X_1 = X_1, ..., X_t = X_t) = \frac{\frac{(\frac{\epsilon}{1-\epsilon})^{\sum_{s=1}^{t} x_s} (1-\epsilon)^{t+1} + (\frac{1-\epsilon}{\epsilon})^{\sum_{s=1}^{t} x_s} \epsilon^{t+1}}{(\frac{\epsilon}{1-\epsilon})^{\sum_{s=1}^{t} x_s} (1-\epsilon)^{t} + (\frac{1-\epsilon}{\epsilon})^{\sum_{s=1}^{t} x_s} \epsilon^{t}}}$$
(2)

## 深度解读

这个数学例子是全文论证的核心,它将前面"天体运行"的比喻用严格的数学语言表达了出来。让 我们来拆解这个例子: 想象一个游戏:

- 1. 我(上帝)先偷偷抛一枚"魔法硬币" Z。它有一半概率是正面(Z=1),一半概率是反面(Z=0)。这个结果我藏起来不告诉你。
- 2. 然后,我抛 T 次 "普通硬币"  $X_1, X_2, ..., X_T$ 。这些普通硬币是有偏的:它们绝大多数情况下会和我那枚魔法硬币的结果一样,但有很小的概率  $\epsilon$  会 "出错"(比如魔法硬币是正面,但普通硬币却显示了反面)。

现在,你的任务是预测下一次普通硬币 $X_{t+1}$ 的结果。

情况一: 你知道魔法硬币Z的结果(包含潜变量的模型) 如果你知道魔法硬币是正面(Z=1),那么你知道下一次普通硬币有  $1-\epsilon$  的概率也是正面。如果魔法硬币是反面(Z=0),那么下一次普通硬币有  $\epsilon$  的概率是正面。这就是公式(1)所表达的。它非常简洁明了。

情况二:你不知道魔法硬币Z的结果,只能看已经抛出的普通硬币序列(纯自回归模型) 这时,事情变得极其复杂。你每看到一次普通硬币的结果,都必须更新你对"那枚看不见的魔法硬币究竟是正面还是反面"的猜测。

- 如果你看到了一连串的正面,你会越来越相信魔法硬币是正面的。
- 如果你看到了一连串的反面,你会越来越相信魔法硬币是反面的。
- 如果你看到的正反面交替出现,你的计算会更加纠结。

公式(2)就是这个复杂"猜测+更新"过程的数学表达。你看它有多么复杂!它需要计算到目前为止所有正面出现的次数( $\sum X_s$ ),并带入一长串复杂的指数和分数运算。

这两个公式的巨大反差,直观地证明了作者的观点: 当一个系统天然存在一个核心的潜变量 Z 时,一个知道 Z 的模型(公式1)会比一个不知道 Z、只能从观测数据 X 中反向推断的模型(公式2)要简单得多。强行使用自回归模型,就是逼迫AI去执行公式(2)这样复杂且低效的运算。

### 原文翻译

We could easily come with worse examples when expressed autoregressively, for instance when the latent variables are positions in the sequence, e.g. the index where a certain pattern occurs as in the example of § 4.1. What we observe in such cases is that it requires running estimates of probabilities ("probability that the target appears here") for which estimation errors are unavoidable and problematic.

当用自回归方式表达时,我们很容易想出更糟糕的例子,例如当潜变量是序列中的位置时,比如在 § 4.1的例子中某个特定模式出现的索引。在这种情况下,我们观察到模型需要对概率进行动态估计 (例如,"目标出现在此处的概率"),而这种估计的误差是不可避免且会引发问题的。

## 深度解读

作者在这里进一步扩展了潜变量的应用场景。潜变量不仅可以是"类别"(正面/负面),还可以是"位置"。想象一下,让AI生成一个故事,要求故事的"转折点"必须出现在第500个词左右。对于标准模型,它必须一边写,一边"惦记"着:"我现在写到第几个词了?是不是快到500了?我是不是该准备转折了?"这种"动态估计"非常困难,很容易出错。而对于Free Transformer,它可以先设定一个潜变量 Z= "转折点位置在第500词",然后从容地生成前面的铺垫,并在接近500词时自然地引入转折。这种"先规划,后执行"的模式显然更优越。

## 原文翻译

The consequence is that a purely auto-regressive density model suffers potentially from several drawbacks:

- It requires an unnecessarily complicated computation, and greater capacity, to implicitly make post-hoc decisions or infer latent quantities from the generated tokens.
- It may be sent off track during the process if, by mistake, a few tokens generated are erroneous, ambiguous or contradictory with those generated previously.
- Key concepts do not appear spontaneously due to the "natural" factorization of the distribution, but are built post-hoc by necessity to fit the training samples better. This may be a fundamental weakness when operating out of distribution.

因此,一个纯自回归的密度模型可能遭受以下几个缺点:

- 它需要不必要的复杂计算和更大的模型容量,来从已生成的词元中间接地做出事后决策或推断潜在量。
- 如果在生成过程中,由于错误,生成了几个有误、模棱两可或与之前生成内容相矛盾的词元, 模型可能会偏离轨道。
- 关键概念不是因为分布的"自然"分解而自发出现,而是为了更好地拟合训练样本而事后构建的。当模型在分布外(out of distribution)操作时,这可能是一个根本性的弱点。

## 深度解读

这里总结了纯自回归模型的三个核心弊病,也是本文要解决的关键问题:

- 1. **高昂的隐性成本**:为了模拟潜在的结构,模型被迫变得更复杂、参数更多("更大的模型容量"),这就像是为了实现一个简单的功能而写了一堆"垃圾代码",效率低下。
- 2. **脆弱性/易偏离**:模型的决策完全依赖于历史生成,就像一个走钢丝的人,一步走错就可能满盘皆输。一个早期的错误词语可能会让整个生成过程走向一个完全错误的方向,并且难以纠正。
- 3. **泛化能力差**:模型学到的"概念"(比如正面/负面)是它为了拟合数据而"硬凑"出来的,而不是从数据的内在结构中"自然"学到的。这导致当遇到训练数据中没见过的新情况时("分布外操作"),模型可能会表现得很差,因为它没有学到真正通用的、底层的规律。

### 原文翻译

The main objective of the present work is to address these issues by providing the model with the freedom of conditioning its auto-regressive process on latent random quantities that are not imposed by the training examples. For instance, for the review generator example above, the model could use a random Boolean value to decide once for all whether the tokens it produces are from the distribution of negative or positive reviews, removing the need for a complicated posterior estimate from the tokens already generated.

本工作的主要目标是通过赋予模型一种自由——让其自回归过程能够以潜在的随机量为条件,而这些随机量并非由训练样本强加——来解决这些问题。例如,在上面提到的评论生成器例子中,模型可以使用一个随机的布尔值来一劳永逸地决定它所产生的词元是来自负面评论的分布还是正面评论的分布,从而无需再根据已生成的词元进行复杂的后验估计。

## 深度解读

这一段清晰地阐述了本文的解决方案和其核心思想——"自由"。论文标题"The Free Transformer"的"Free"就来源于此。这种"自由"体现在模型可以自己创造和使用这些潜变量,而不是被动地从数据中推断。作者用"一劳永逸"(once for all)这个词来形容这种决策方式的优越性。一旦决定了评论的基调,模型就可以心无旁骛地进行创作,而不需要在每一步都瞻前顾后、反复进行"复杂的后验估计"。这正是从"战术家"到"战略家"的转变。

# 3. 方法

## 原文翻译

Any latent random value  $Y_r$  whatever its statistical dependency with the tokens  $S_1$ , ...,  $S_l$  and other latent  $Y_1$ , ...,  $Y_{r-1}$  sampled so far, can be expressed under reasonable assumptions as  $f_r(S_1, ..., S_t, Y_1, ..., Y_{r-1}, Z_r)$  where  $Z_r$  is a value coming from a random generator. Hence, if we provide the model with enough random values  $Z_1$ ,  $Z_2$ , ... sampled independently during generation, a proper training procedure could in principle build families of latent variables with arbitrary dependency structure, as long as the model's capacity allows it to encode  $f_r$ .

在合理的假设下,任何潜在随机值  $Y_r$ ,无论其与目前已采样的词元  $S_1$ ,…,  $S_l$  和其他潜变量  $Y_1$ ,…,  $Y_{r-1}$  的统计依赖关系如何,都可以表示为  $f_r(S_1,...,S_t,Y_1,...,Y_{r-1},Z_r)$ ,其中  $Z_r$  是一个来自随机生成器的值。因此,如果我们在生成过程中为模型提供足够多的独立采样的随机值  $Z_1,Z_2$ ,…,只要模型的容量允许它编码函数  $f_r$ ,一个合适的训练程序原则上可以构建出具有任意依赖结构的潜变量族。

## 深度解读

这一段从理论上阐述了方法的可行性。其核心思想是,任何复杂的决策(潜在随机值  $Y_r$ )都可以被看作是一个函数  $f_r$  的输出,这个函数的输入包括已有的信息(历史词元和其他决策)和一个纯粹的随机噪声源  $Z_r$ 。这就像一个厨师做菜,最终菜品的味道( $Y_r$ )取决于已有的食材(S 和其他 Y)以及厨师的临场发挥(随机源  $Z_r$ )。作者的论点是,只要我们能给模型一个强大的函数逼近器(即神经网络本身)和源源不断的"随机性"来源(Z),理论上模型就能学会做出任何复杂的、有结构的潜在决策。

## 原文翻译

In the same way that the choice of a token during sampling can be expressed as a function of a random value and the logits, any activation which is a function of a random value and other activations can be interpreted as a decision made by the model during the generative process. Such decisions make the latent activation non-deterministic functions of the tokens, and observing the latter only gives a partial information about the former.

就像采样过程中选择一个词元可以表示为随机值和logits的函数一样,任何一个作为随机值和其他激活值的函数的激活值,都可以被解释为模型在生成过程中做出的一个决策。这样的决策使得潜在激活值成为词元的非确定性函数,观察后者只能提供关于前者的部分信息。

## 深度解读

这里将"潜变量"的概念从一个抽象的外部变量,具体化为模型内部的"激活值"(activation)。在神经网络中,每一层的输出都是激活值。作者提出,如果我们在网络中间的某个环节,将一个随机噪声注入到激活值的计算中,那么这个激活值本身就变成了一个"潜在决策"。因为它的值不仅取决于输入的文本,还取决于这个随机噪声,所以它变得"非确定性"了。即使输入完全相同的文本,每次生成的内部"潜在决策"也可能不同。这就为模型提供了"自由"的源泉。观察最终生成的文本,只能部分地反推出模型当时内心可能做出了什么样的"潜在决策"。

### 图1描述

- 图(a) 展示了一个标准的解码器Transformer。数据流从底部的输入  $S_{1:T-1}$  开始,经过解码器(Dec),最终在顶部输出对下一个词元的预测  $\hat{P}(S_{2:T})$ 。这是一个单向的、确定性的过程。
- 图(b) 展示了如何将标准模型扩展为一个条件变分自编码器(cVAE)。在推理(生成)时(左图),模型接收一个随机采样的潜变量 Z和输入  $S_{1:T-1}$ ,共同生成输出。在训练时(右图),模型需要一个额外的编码器(Enc),它读取整个序列  $S_{1:T-1}$  来生成一个与该序列匹配的 Z,然后解码器再利用这个 Z来重构序列。这种方式的问题是,编码器和解码器是两个独立的、庞大的模型,成本很高。

• 图(c) 展示了本文提出的"自由Transformer"(Free Transformer)的巧妙设计。它将解码器分成了两半(Dec 1/2 和 Dec 2/2)。在推理时(左图),随机的 Z 在解码器的中间层被注入。在训练时(右图),编码器(Enc)巧妙地复用了第一半解码器(Dec 1/2),只增加了一个专门的、非因果的编码层。这样一来,编码器的绝大部分计算和参数都与解码器共享了,极大地降低了额外开销。这正是该方法在工程上如此高效的关键所在。

## 3.1 条件变分自编码器

## 原文翻译

Generating a full sequence from scratch with a model that depends on a random variable Z is trivial: sample  $Z \sim P(Z)$  and then run the standard auto-regressive process, with the computation of the logits modulated by Z.

使用一个依赖于随机变量Z的模型从头生成一个完整序列是很简单:从  $Z \sim P(Z)$  中采样一个Z,然后运行标准的自回归过程,其中logits的计算会受到Z的调节。

## 深度解读

这里描述的是模型的"使用"阶段(即推理或生成)。这个过程非常直观。就像我们之前比喻的作家:

- 1. 首先,从"所有可能的主题"中随机挑选一个主题(采样 Z)。比如,模型随机采样了一个代表"科技新闻风格"的潜变量 Z。
- 2. 然后,模型就像一个普通的GPT一样开始逐字生成,但它的每一步决策都会受到这个 Z 的影响。它会更倾向于使用"芯片"、"人工智能"、"市场"等词汇,并采用相应文体的句式。这个过程本身仍然是自回归的,但它被一个全局的"条件"所引导。

## 原文翻译

Training the model, however, is far more involved. Given a training sample S, the objective is to maximize  $P(S) = \int_{Z} P(S \mid Z = z) P(Z = z) dz$ , (3) which can be estimated only if we can get Zs consistent with S, which amounts to a complex inference problem if we want Z to capture meaningful structural properties of the sequence.

然而,训练模型要复杂得多。给定一个训练样本S,目标是最大化  $P(S) = \int_{Z} P(S \mid Z = Z) P(Z = Z) dZ$ ,(3) 只有当我们能得到与S一致的Z时,这个值才能被估计。如果我们希望Z能捕捉序列有意义的结构特性,这就构成了一个复杂的推断问题。

### 深度解读

模型的"训练"阶段则要困难得多。我们的目标是让模型在看到大量文本后,学会生成类似的文本。公式(3)表达了这个目标:我们希望模型生成样本S的总概率P(S)最大化。这个总概率是所有可能的"隐藏计划"Z下生成S的概率 $P(S \mid Z = Z)$ 的加权平均。这里的难题在于,对于一篇给定的文章S(比如一篇真实的科技新闻),我们并不知道作者当初是基于哪个"隐藏计划"Z写出来

的。我们需要一个方法来"反向推断"出与这篇新闻最匹配的那个 Z。这个"反向推断"问题非常困难,因为可能的 Z有无数个。

## 原文翻译

Providing those Zs is the role of the encoder of a Variational Autoencoder (Kingma and Welling, 2013), whose main purpose is to sample from a "good" distribution  $Q(Z \mid S)$  so that a sampled Z modulates the decoder in a way that leads it to generate S.

提供这些Z,正是变分自编码器(VAE)中编码器的角色(Kingma and Welling, 2013)。其主要目的是从一个"好的"分布  $Q(Z \mid S)$  中采样,使得采样出的Z能够调节解码器,引导它去生成S。

## 深度解读

这里引入了解决上述难题的工具:编码器(Encoder)。让我们用一个"艺术家-评论家"的类比来理解VAE的训练过程:

- **解码器(Decoder)**:一个有才华但缺乏灵感的"艺术家"。你给他一个抽象的灵感(潜变量 Z),他就能画出一幅画(生成文本 S)。
- **编码器(Encoder)**:一个眼光毒辣的"艺术评论家"。你给他一幅现成的名画(训练文本 S),他能精准地提炼出这幅画的核心"灵魂"或"灵感"(推断出对应的 Z)。这个过程就是  $Q(Z \mid S)$ 。

训练过程就像是让艺术家和评论家进行配对练习:

- 1. 拿一幅名画 S 给评论家看。
- 2. 评论家提炼出他认为的"灵感"Z。
- 3. 把这个"灵感" Z交给艺术家。
- 4. 艺术家根据这个 Z 重新创作一幅画。
- 5. 我们希望艺术家重画出来的画,和最初那幅名画越像越好。

通过成千上万次的这种练习,评论家(编码器)的反向推断能力和艺术家(解码器)的正向创作能力都得到了提升。

### 原文翻译

We follow this approach and optimize jointly the parameters of the decoder and the parameters of a second model, which is an encoder in the VAE sense. Even though the noise Z has no relation to S initially, if the training succeeds, the model will use it to structure the generative process. In the example of a movie review generator of the previous section, for instance, given a review from the training set, the encoder would implicitly classify it as positive or negative, and generate a consistent Z. Increasing  $P(S \mid Z)$  with that Z could be interpreted as improving the "negative review generator" or the "positive review generator" that are implicitly encoded in the decoder's weights.

我们遵循这种方法,并联合优化解码器的参数和第二个模型(即VAE意义上的编码器)的参数。尽管噪声Z最初与S没有任何关系,但如果训练成功,模型将使用它来结构化生成过程。例如,在上一节的电影评论生成器例子中,给定一个来自训练集的评论,编码器会隐式地将其分类为正面或负面,并生成一个与之一致的Z。用这个Z来增加  $P(S \mid Z)$ ,可以被解释为在改进解码器权重中隐式编码的"负面评论生成器"或"正面评论生成器"。

## 深度解读

这一段解释了训练成功后的效果。一开始,潜变量 Z 只是一个随机噪声,没有任何意义。但通过上述的"艺术家-评论家"训练,模型会自发地赋予 Z 的不同取值不同的含义。例如,编码器在看到大量负面评论后,会学会将它们都映射到潜空间中的某个特定区域(比如,输出一个特定的 Z 值或 Z 的分布)。然后,解码器在接收到这个区域的 Z 时,就会知道应该激活"负面评论生成"模式。这样,模型就"无监督"地学会了"正面/负面"这个概念,并将其与潜变量 Z 绑定。训练过程中的"提高  $P(S \mid Z)$ ",实际上就是在不断打磨和完善这些被 Z 索引的、各种各样的"子生成器"。

## 原文翻译

A key element of this approach is to limit the amount of information flowing from the encoder to the decoder through Z, so that the encoder does not provide quantities that should be computed by the decoder. At the limit the encoder could copy entirely S into Z so that a trivial decoder, useless without the encoder, hence in inference, would score perfectly in training. The formal derivation of the VAE shows that the proper measure of information is the Kullback-Leibler divergence between  $Q(Z \mid S)$  and P(Z), and that the loss to minimize should sum it with the reconstruction loss.

这种方法的一个关键要素是限制通过Z从编码器流向解码器的信息量,以防止编码器提供本应由解码器计算的量。在极限情况下,编码器可以将S完全复制到Z中,这样一来,一个平凡的解码器(在没有编码器的情况下毫无用处,因此在推理时也无用)在训练中会得到完美的分数。VAE的正式推导表明,衡量信息量的恰当指标是  $Q(Z \mid S)$  和 P(Z) 之间的KL散度(Kullback-Leibler divergence),并且要最小化的损失函数应该将它与重构损失相加。

#### 深度解读

这是VAE训练中的一个核心要点,也是一个非常巧妙的设计,被称为"信息瓶颈"(Information Bottleneck)。 想象一下,如果不对评论家(编码器)提炼的"灵感" Z做任何限制,他可能会"作弊"。为了让艺术家能完美复原名画,他可能会写一份极其详尽的笔记,把名画的每一个像素、每一笔触都记录下来,然后把这份笔记交给艺术家。艺术家照着抄,当然能画出一模一样的画。但这样训练出来的艺术家并没有学会真正的创作,他只是一个复印机。在实际使用中(推理阶段),我们没有名画给评论家看,也就无法产生详尽的笔记,艺术家就什么也画不出来了。 为了防止这种情况,我们必须给评论家提一个要求:你的"灵感"笔记 Z必须非常简洁,并且要符合某种简单的、预设的格式(比如,必须像一个标准正态分布 P(Z))。KL散度就是衡量评论家笔记的实际格式  $Q(Z \mid S)$  与我们要求的标准格式 P(Z) 之间差异的数学工具。 因此,总的训练目标(损失函数)就包含两部分:

1. 重构损失:要求艺术家画的画要像原作。

2. **KL散度损失**:要求评论家的笔记必须简洁且符合标准格式。这两者是相互制约的。评论家为了让艺术家画得更像,会倾向于在笔记里塞入更多信息,导致KL散度增大;而为了降低KL散度,他又必须压缩信息,只保留最核心、最精华的部分。这种制衡迫使编码器学会真正的"抽象"和"概括",而不是简单的"复制粘贴",从而让解码器学到真正的生成能力。

### 3.2 模型结构

## 原文翻译

In what follows, we call "Transformer Block" the usual combination of a Multi-Head Attention layer and a MLP-like tokenwise module, with normalisation layers and residual connections.

在下文中,我们将"Transformer模块"(Transformer Block)称为多头注意力层(Multi-Head Attention)和类MLP的逐词元模块(tokenwise module)的常规组合,并包含归一化层和残差连接。

### 深度解读

这是一个术语定义。作者明确了"Transformer模块"的具体构成,它包含了两个核心部件:多头注意力机制(负责捕捉词与词之间的关联)和前馈神经网络(MLP,负责对每个词的信息进行独立加工),以及一些帮助模型稳定训练的标准化组件。一个完整的Transformer模型就是由许多这样的模块堆叠而成的。

## 原文翻译

As pictured on Figure 1 and 2, the Free Transformer is a standard decoder with a noise Z injected in its middle layer. This allows to share half of the Transformer block with the encoder, cutting down drastically the computational overhead by having a single Transformer block that has to be computed specifically for the encoder. Hence, as we will see, this model possesses all the components of a decoder Transformer and has an additional non-causal block and two linear layers for the encoder. While we did not investigate what is the best depth to inject Z, doing it too early would reduce the encoder's capacity, and doing it too late would reduce the decoder's capacity to process the latent variables.

如图1和图2所示,自由Transformer是一个标准的解码器,其噪声Z被注入到其中间层。这使得编码器可以与解码器共享一半的Transformer模块,通过只让一个Transformer模块专用于编码器计算,从而大幅削减了计算开销。因此,正如我们将看到的,该模型拥有解码器Transformer的所有组件,并为编码器增加了一个额外的非因果模块和两个线性层。虽然我们没有研究注入Z的最佳深度,但注入得太早会降低编码器的能力,而注入得太晚则会降低解码器处理潜变量的能力。

### 深度解读

这里揭示了Free Transformer架构设计的核心巧思,也是其能够保持"低开销"的关键。传统的VAE模型,编码器和解码器是两个独立的网络,总参数量约等于两者之和。而Free Transformer的设计

是"融合"而非"拼接"。它让编码器和解码器"共享"了前半部分网络。 具体来说,输入文本首先流经前一半的Transformer模块,得到一个中间表示。然后:

- 编码器路径:这个中间表示被送入一个专用的、小型的编码器模块,用于推断潜变量 Z。
- **解码器路径**: 这个中间表示与根据 *Z* 生成的信息相融合,然后继续流经后一半的Transformer 模块,最终生成输出。

这种设计的好处是显而易见的:编码器所需的额外参数和计算量非常小,因为它的大部分"身体"都借用了解码器的。作者也提到了注入位置的权衡:注入点太靠前,编码器能利用的信息太少,难以做出好的推断;注入点太靠后,解码器没有足够的层次来消化和利用 Z 带来的信息。选择中间位置是一个合理的折中。

### 原文翻译

For clarity, we omit in what follows the batch size in the tensor shapes. As a standard decoder Transformer, the Free Transformer processes a sequence of tokens by first encoding them with the embedding table into a tensor  $X_0$  of shape  $T \times D$ . Then it evaluates sequentially the first L/2 Transformer blocks to get  $X_{L/2}$  of same shape, and at this point, it samples a sequence of one-hot vectors  $Z = (Z_1, ..., Z_t) \in \{0, 1\}^{T \times C}$ .

为清晰起见,我们在下文的张量形状中省略了批量大小(batch size)。作为一个标准的解码器 Transformer,自由Transformer通过首先使用嵌入表(embedding table)将词元序列编码为一个 形状为  $T \times D$  的张量  $X_0$  来处理它。然后,它顺序评估前 L/2 个Transformer模块,得到同样形状的  $X_{L/2}$ ,此时,它采样一个独热(one-hot)向量序列  $Z = (Z_1, ..., Z_t) \in \{0, 1\}^{T \times C}$ 。

## 深度解读

这里开始描述数据在模型中流动的具体过程。

- 1. **嵌入(Embedding)**: 输入的文本(一串词元)首先被转换成计算机能够理解的数学形式——向量。每个词元都对应一个高维向量。这样,一个长度为 T 的句子就变成了一个  $T \times D$  的矩阵(或称张量),其中 D 是向量的维度。
- 2. **前向传播(前半部分)**: 这个矩阵依次通过前 L/2 层的Transformer模块,每一层都会对它进行加工和提炼,得到中间表示  $X_{L/2}$ 。
- 3. **采样Z**:在网络的中点,模型会为序列中的每一个位置 t 采样一个潜变量  $Z_t$ 。这里的  $Z_t$  是一个"独热向量",即一个长度为 C、只有一个元素为1其余都为0的向量。你可以把它想象成一个有 C 个选项的开关,模型在每个时间步都选择打开其中一个。

## 原文翻译

During generation, this is done by sampling, for each  $Z_t$ , an index c uniformly in  $\{0, ..., C-1\}$ , and then encoding it as a one-hot vector of dimension C. During training or KV cache prefilling, Z has to be consistent with the tokens of S already fixed, and the sampling is done with the encoder instead, as described in § 3.3.

在生成阶段,这是通过为每个  $Z_t$  在  $\{0,...,C-1\}$  中均匀采样一个索引c,然后将其编码为维度为C 的独热向量来完成的。在训练或KV缓存预填充阶段,Z必须与S中已固定的词元保持一致,此时采样将由编码器完成,详见  $\S$  3.3。

### 深度解读

这里区分了两种模式下Z的来源:

- **生成(推理)模式**: 此时我们没有可供参考的完整文本,模型需要"自由创作"。因此,它直接随机地、均匀地选择一个潜变量  $Z_t$ 。这就像作家在构思时,随机从脑海中的无数个灵感中挑选一个。
- 训练(或预填充)模式:此时我们有一篇完整的文章 S 作为学习材料。我们需要推断出与这篇文章最匹配的 Z。这个工作就交给了我们之前提到的"评论家"——编码器。编码器会读取整篇文章,然后为每个位置生成一个最合适的  $Z_t$ 。

## 原文翻译

This tensor Z is processed by a linear layer to obtain a tensor R of shape  $T \times D$ . Then, the L/2+1-th Transformer block gets as input for queries the tensor  $X_{L/2}$  and as input for keys and values the tensor  $X_{L/2}+R$ . The rest of the Transformer blocks are evaluated in sequence to get  $X_L$  which is processed by the read-out linear layer to obtain the logit tensor L of shape  $T \times V$ , where V is the vocabulary size.

这个张量Z经过一个线性层处理,得到一个形状为  $T \times D$  的张量R。然后,第 L/2+1 个 Transformer模块接收张量  $X_{L/2}$  作为查询(queries)的输入,并接收张量  $X_{L/2}+R$  作为键(keys)和值(values)的输入。剩余的Transformer模块被顺序评估,得到  $X_L$ ,它再经过读出线性层处理,得到形状为  $T \times V$  的logit张量L,其中V是词汇表大小。

## 深度解读

这是潜变量 Z如何与主干信息流融合的关键步骤,也是一个非常重要的架构细节。

- 1. 采样得到的 Z(一个  $T \times C$  的矩阵)首先经过一个线性层(可以看作一个简单的神经网络层)变换,变成一个与主干信息流  $X_{L/2}$  形状完全相同的张量 R ( $T \times D$ )。
- 2. 然后,在第 L/2+1 层的注意力机制中,发生了一个巧妙的操作:查询(Query)仍然来自原始的中间表示  $X_{L/2}$ ,但键(Key)和值(Value)则变成了  $X_{L/2}+R$ 。这个"加法"操作意义非凡。它不是简单地把两个信息拼接在一起,而是用 R作为一个"引导向量"或"偏置项",对原始的文本表示  $X_{L/2}$  进行了整体的、方向性的"微调"。如果 Z编码了"负面情绪",那么变换后的 R就会将  $X_{L/2}$  中的所有词向量都朝着"负面"的语义方向推动一下。这样,后续的Transformer模块在进行信息交互时,就会自然而然地更关注与负面情绪相关的信息,并最终生成负面词汇。这是一种非常高效和优雅的全局信息注入方式。
- 3. 经过融合后,数据流继续通过剩下的 L/2 个模块,最后通过一个"读出层"将最终的表示  $X_L$  转换成对词汇表中每个词的预测概率(logits)。

### 图2描述

这张图是Free Transformer架构的"工程蓝图",比图1更加详细。让我们从下往上,一步步追踪数据的旅程:

- 1. **输入(S1:T-1)**: 文本序列进入模型。
- 2. 嵌入层 (Embeddings): 文本被转换为  $T \times D$  的向量矩阵。
- 3. **解码器前半部分 (Decoder Causal Transformer Block x L/2)**:数据流经前 L/2 个标准的、带 因果掩码的Transformer模块。这是编码器和解码器共享的部分。
- 4. 分岔点: 在网络的中间,数据流分岔。
  - 编码器路径(橙色部分,仅在训练时激活):
    - 前半部分的输出  $X_{L/2}$  作为键(KV)和值,一个可学习的特殊向量  $\xi$  作为查询(Q),共同输入到一个**非因果**的Transformer模块(Encoder Non-Causal Transformer Block)。 "非因果"意味着在这一步,为了推断Z,模型可以同时看到整个句子,不受从左到右顺序的限制。
    - 该模块的输出经过一个线性层(Encoder read-out FC),得到维度为  $T \times H$  的输出。
    - 这个输出被送入**二进制映射器 (Binary Mapper)**,最终采样得到潜变量 Z (形状为 $T \times 2^H$ )。

## • 解码器路径:

- 在推理时, Z 是通过均匀采样得到的。在训练时, Z 来自编码器路径。
- Z经过一个线性层(Post-sampler FC),得到引导向量 R。
- 在第 L/2 + 1 个模块(Decoder Causal Transformer Block)中,前半部分的输出  $X_{L/2}$  作为Q,而  $X_{L/2} + R$ 作为KV,实现了信息的融合。
- 5. **解码器后半部分 (Decoder Causal Transformer Block x L/2-1)**: 融合后的信息流经剩下的 Transformer模块。
- 6. 输出层 (Decoder read-out FC): 最终的表示被转换为对词汇表的预测 logits。

图中橙色部分是Free Transformer比标准Transformer多出来的组件,它们只在训练时使用,因此推理时的额外开销几乎为零。带虚线框的算子(如Binary Mapper)表示它们没有可训练的参数。这个设计体现了极高的工程效率。

### 3.3 编码器和损失函数

#### 原文翻译

As stated in the previous section, during training or KV cache pre-filling, the tensor Z is sampled with the encoder. The Free Transformer possesses one Transformer block specific to the encoder, which is non-causal, making the encoder as a whole non-causal. This is necessary since the conditioning by the decoder may have long-range effects, requiring the full sequence to be taken into account to get a proper conditional distribution of the latent.

如前一节所述,在训练或KV缓存预填充期间,张量Z是用编码器采样的。自由Transformer拥有一个专用于编码器的Transformer模块,该模块是非因果的,从而使整个编码器成为非因果的。这是必要的,因为解码器的条件作用可能具有长程效应,需要考虑整个序列才能获得潜变量的正确条件分布。

## 深度解读

这里再次强调了编码器模块的"非因果"(non-causal)特性。标准的解码器是"因果"的,即在预测第t个词时,它只能看到前t-1个词。但编码器的任务是为整个句子S提炼一个"总结"Z,它必须能看到句子的全貌才能做出准确的判断。比如,要判断一个句子的情感,你通常需要读完整句话。因此,这个专用的编码器模块必须打破从左到右的顺序限制,允许信息在整个序列中自由流动。

## 原文翻译

This encoder-specific block gets as input for the queries a trained token embedding  $\xi$  replicated to match the sequence length, and for the keys and values the output of the first half of the decoder's blocks. The motivation for using a learned constant input for the queries instead of the standard representation of the input sequence is to prevent the encoder from building a token-wise mapping and make it instead capture global properties of the sequence that may be more transferable across tasks and data-sets.

这个编码器专用模块接收一个经过训练的词元嵌入  $\xi$ (复制以匹配序列长度)作为查询的输入,并接收解码器前半部分模块的输出作为键和值的输入。使用一个可学习的常量作为查询输入,而不是使用输入序列的标准表示,其动机是防止编码器建立逐词元的映射,而是促使它捕捉序列的全局属性,这些属性可能在不同任务和数据集之间更具可迁移性。

## 深度解读

这是一个非常精妙的设计选择。在标准的注意力机制中,Q、K、V通常都来源于同一份输入。但在这里,编码器的Q是一个特殊的、可学习的常量向量  $\xi$ 。这意味着编码器在"提问"时,问的总是同一个"全局性的问题",比如"请总结一下这个序列的整体风格是什么?"或者"这个序列的核心主题是什么?"。它不是在问"第i个词是什么意思?"。这种设计强迫编码器去关注和捕捉整个序列的"全局属性",而不是陷入对局部细节的逐词元分析。这样学到的潜变量Z更有可能代表高层次的、抽象的概念,从而在不同的任务中都具有泛化能力。

### 原文翻译

A linear readout computes from the encoder block's output a vector of dimension H=16 for every token.

一个线性读出层从编码器模块的输出中为每个词元计算一个维度为 H=16 的向量。

## 深度解读

编码器模块处理完信息后,一个简单的线性层会将其输出转换为一个维度为16的向量。这个向量将作为下一步生成潜变量 Z的基础。

## 算法1和算法2的描述

## 算法1:标准解码器Transformer的前向传播

- 1. FORWARD(tokens):定义函数,输入为词元序列。
- 2. x ← embeddings(tokens):将词元转换为向量。
- 3. for n = 1,..., B do:循环B次(B为总层数)。
- 4. x ← blocks[n](in=x):将x依次送入每个Transformer模块。
- 5. end for:结束循环。
- 6. logits ← linear\_readout(RMS\_norm(x)):对最终输出进行归一化和线性变换,得到预测logits。
- 7. return logits:返回logits。这是一个简单、线性的流程。

### 算法2: 自由Transformer的前向传播

- 1. FORWARD(tokens):定义函数。
- 2. x ← embeddings(tokens):词元转向量。
- 3. for n = 1, ..., B/2 do: 只循环前B/2次。
- 4. x ← blocks[n](in=x):数据流经前半部分网络。
- 5. end for:前半部分结束。
- 6. if train or prefill then:判断是否为训练模式。
- 7. y ← encoder block(in g = zeta, in kv = x):如果是训练,则激活编码器模块(zeta即\$\xi\$)。
- 8. o ← encoder linear readout(RMS norm(y)):得到编码器输出。
- 9. z ← binary mapper(o):通过二进制映射器采样z。
- 10. else : **如果是生成模式**。
- 11. z ← one hot(uniform sampler()):随机均匀采样z。
- 12. end if:模式判断结束。
- 13. r ← linear post sampler(z):将z变换为引导向量r。
- 14. x ← blocks(in q=x, in kv=x+r):在中间层将r和x融合。
- 15. for n = B/2+2,..., B do:循环后半部分的层。
- 16. x ← blocks[n](in=x):数据流经后半部分网络。
- 17. end for:结束循环。
- 18. logits ← linear readout(RMS norm(x)):得到最终logits。
- 19. return logits:返回。算法2清晰地展示了在网络中间发生的"分岔-融合"过程,以及训练和生成模式的差异。

### 原文翻译

These components are interpreted as logits of individual bit, used to sample a value in  $\{0,...,2^H-1\}$  which is encoded into a one-hot vector of dimension  $2^H=65,536$ , with gradient pass-through, as described in § 3.4. Hence, the random embedding Z is a sequence of T one-hot vectors  $Z_t$  of dimension  $2^H$ . The prior distribution used for generation is uniform  $P(Z_t=z)=1/2^H$ , and  $Q(Z\mid S=s)$  is the distribution corresponding to the sampling with the encoder described above.

这些分量被解释为单个比特的logits,用于在  $\{0,...,2^H-1\}$  中采样一个值,该值被编码为维度  $2^H=65,536$  的独热向量,并带有梯度直通,如§ 3.4所述。因此,随机嵌入Z是一个由T个维度为  $2^H$  的独热向量  $Z_t$  组成的序列。用于生成的先验分布是均匀分布  $P(Z_t=z)=1/2^H$ ,而  $Q(Z\mid S=s)$  是与上述使用编码器采样相对应的分布。

## 深度解读

这里解释了如何从那个 H=16 维的向量生成最终的潜变量 Z。模型不直接从65536个选项中选一个(这在计算上很难处理),而是把它分解成16个独立的"是/否"决策(二进制位)。这16个二进制位组合起来,就可以表示从0到65535的任何一个数字。最终,这个数字被转换成一个65536维的独热向量,作为潜变量  $Z_t$ 。在生成时,我们假设这65536个可能性是等概率的(均匀先验分布)。在训练时,编码器给出的概率分布  $Q(Z \mid S)$  则是由这16个二进制位的概率决定的。

## 原文翻译

深度解读

The KL divergence is then equal to  $D_{KL}(Q(Z_t \mid S_1,...,S_T) \mid P(Z_t)) = H \log 2 + \sum_{z=1}^{2^H} Q(Z=z \mid S) \log Q(Z=z \mid S)$ . (4) We control it by adding it to the loss, and prevent its collapse by using a token-wise free bits method (Kingma et al., 2016). This means that we sum the KL divergence of individual  $Z_t$  that are above a threshold and ignore the others. This leads us to use for training loss the sum of the standard cross-entropy and the following quantity  $\frac{1}{T} \sum_{t=1}^{T} \max(0, D_{KL}(Q(Z_t \mid S_1,...,S_T) \mid P(Z_t)) - \kappa)$  (5) where the threshold  $\kappa$  is an hyperparameter.

那么KL散度等于  $D_{KL}(Q(Z_t \mid S_1,...,S_T) \mid P(Z_t)) = H \log 2 + \sum_{z=1}^{2^H} Q(Z = z \mid S) \log Q(Z = z \mid S)$ . (4) 我们通过将其加入损失函数来控制它,并使用逐词元的"自由比特" (free bits) 方法(Kingma et al., 2016)来防止其坍缩。这意味着我们只对那些超过阈值的单个  $Z_t$  的KL散度求和,而忽略其他的。这使我们在训练时使用的损失函数是标准交叉熵与以下量的总和  $\frac{1}{T} \sum_{t=1}^{T} \max(0, D_{KL}(Q(Z_t \mid S_1,...,S_T) \mid P(Z_t)) - \kappa)$  (5) 其中阈值  $\kappa$ 是一个超参数。

这里详细说明了损失函数的设计,特别是如何处理KL散度。 公式(4)是计算KL散度的具体公式,当先验是均匀分布时,它可以简化为熵的形式。 核心是公式(5)引入的 "自由比特"(free bits)技巧。前面我们提到,KL散度是为了防止编码器 "作弊",给它一个 "信息瓶颈"。但如果这个瓶颈太窄,编码器可能会觉得 "多一事不如少一事",干脆完全不使用潜变量Z,让KL散度一直为0。这种情况被称为 "后验坍塌"(posterior collapse),模型退化成了一个标准的自回归模型。 "自由比特"就像是给编码器一个 "信息预算"  $\kappa$ 。公式(5)的含义是:

- 只要你(编码器)使用的信息量(KL散度)没有超过预算 *K*,我就不惩罚你。
- 只有当你使用的信息量超出了预算,我才把超出部分加到总的损失里去惩罚你。 这个简单的 max(0,...) 操作,极大地鼓励了模型去"用满"这个免费的信息额度,从而有效避免了后验坍塌,确保潜变量Z被真正地利用起来。超参数 κ的大小,就决定了我们允许模型在每个时间步的"隐藏计划"中包含多少比特的信息。

## 3.4 二进制映射器

## 原文翻译

The last linear layer of the encoder computes for every index t of the sequence being processed a vector  $L_t = (L_{t,1}, ..., L_{t,H}) \in \mathbb{R}^H$ , whose components are interpreted as the logits of individual bits of a binary encoding. The Binary Mapper samples those bits  $B_{t,1}, ..., B_{t,H}$  independentely with  $P(B_{t,h} = 1) = \frac{1}{1+e^{-L_{t,h}}}$  (6) and outputs a one-hot vector  $Y_t$  of dimension  $2^H$  corresponding to the resulting value:  $Y_{t,d} = \{1 \text{ if } d = 1 + \sum_{h=1}^H 2^{h-1} B_{h,t} \text{ otherwise.} \}$ 

### 深度解读

这一部分是技术实现细节,解释了"二进制映射器"的工作原理。

- 1. 编码器输出一个16维的向量  $L_t$ 。这16个数字,每一个都代表一个二进制位的"倾向性"(logits)。一个大的正数意味着这个位很可能是1,一个大的负数意味着很可能是0。
- 2. 公式(6)是Sigmoid函数,它将这个"倾向性"数值转换成一个0到1之间的概率。
- 3. 模型根据这个概率,为16个位中的每一个都独立地随机采样一个0或1的值。
- 4. 公式(7)将这16个二进制位组合成一个整数(例如, 100...0 组合成  $2^{15}$ ),然后将这个整数转换成一个65536维的独热向量。这就是最终的潜变量  $Z_t$ (这里用  $Y_t$  表示)。

### 原文翻译

During training, the computation also propagates the gradient of the probabilities of the  $2^H$  values. If  $U(d) = (U_1(d), ..., U_H(d)) \in \{0, 1\}^H$  is the binary encoding of d, and we define  $G_t$  as  $G_{t,d} = P(B_t = U(d-1)) = exp(\sum_h log P(B_{t,h} = U_h(d-1))) = exp(\sum_h (1 - U_h(d-1))) = exp(\sum_h (1$ 

在训练期间,计算过程还会传播这  $2^H$  个值的概率梯度。如果  $U(d) = (U_1(d), ..., U_H(d)) \in \{0,1\}^H$  是d的二进制编码,我们将  $G_t$  定义为  $G_{t,d} = P(B_t = U(d-1)) = \exp(\sum_h \log P(B_{t,h} = U_h(d-1))) = \exp(\sum_h (1 - U_h(d-1)) \log(1 - \frac{1}{1+e^{-L_{t,h}}}) + U_h(d-1) \log(\frac{1}{1+e^{-L_{t,h}}}))$  那么二进制映射器输出  $Y_{t,d} + G_{t,d} - \det C(G_{t,d})$ , (8) 其中  $\det C(X) = X$ 并且其雅可比矩阵  $J_{detach}(X) = 0$ .

## 深度解读

这部分是整个方法中最具技术性的地方,它解决了一个神经网络训练中的普遍难题:如何让梯度"通过"一个随机采样操作。通常,随机采样是一个离散的、不可微分的操作,梯度信息到这里就中断了,导致我们无法用反向传播来训练采样之前的网络层。这里使用了一种被称为"直通估计器"(Straight-Through Estimator)的技巧。公式(8)的含义是:

- 在前向传播时,我们正常进行采样,得到一个离散的独热向量  $Y_{t,d}$ 。
- 在反向传播时,我们"假装"这个操作是可微的。我们计算出每个可能输出的概率  $G_{t,d}$ ,然后把来自后续层的梯度,乘以  $G_{t,d}$  的梯度,传回到前面的网络。  $Y_{t,d}$   $detach(G_{t,d})$  这一项确保了前向传播的值是离散的0或1,而加上  $G_{t,d}$  这一项则在反向传播时提供了一条"平滑"的梯度通路。这是一种工程上的"欺骗",但实践证明它非常有效。

### 原文翻译

The motivation for using a binary encoding instead of having the encoder output  $2^H$  logits directly is to facilitate the gradient pass-through thanks to the monotonicity of the sigmoid.

使用二进制编码,而不是让编码器直接输出  $2^H$  个logits,其动机是利用sigmoid函数的单调性来方便梯度直通。

### 深度解读

最后,作者解释了为什么要费劲地使用二进制编码,而不是让编码器直接预测65536个类别中每个类别的概率。因为预测16个独立的二进制位,每个位的概率都由一个单调的Sigmoid函数控制,这使得梯度计算更加稳定和简单。直接处理65536个类别会引入一个巨大的Softmax层,其梯度行为要复杂得多。这再次体现了该方法在工程实现上的深思熟虑。

## 4. 实验

## 原文翻译

We first test the qualitative behavior of the Free Transformer on a synthetic task in § 4.1, then compare it on multiple benchmarks to baselines with 1.5B and 8B parameters models for various KL divergence thresholds in § 4.4, and finally assess the performance gain of a 8B parameter model trained on 1T tokens in § 4.5.

我们首先在§ 4.1的一个合成任务上测试自由Transformer的定性行为,然后在§ 4.4中,在多个基准上将其与15亿和80亿参数的基线模型在不同KL散度阈值下进行比较,最后在§ 4.5中评估一个在1

万亿(1T)词元上训练的80亿参数模型的性能增益。

### 深度解读

这一段是实验部分的路线图。作者的实验设计非常有条理,层层递进:

- 1. **定性"玩具"实验**: 首先在一个自己设计的、结构已知的简单任务上进行测试。这就像在受控的实验室环境中验证一个物理定律,目的是直观地确认模型是否真的在按照他们设想的方式工作,即利用潜变量Z来编码数据结构。
- 2. **定量对比实验**: 然后在真实、复杂的"下游任务"中,与没有潜变量的标准模型(基线)进行正面比较。实验覆盖了两种主流的模型尺寸(1.5B和8B),并探索了不同的"信息预算" $\kappa$ ( $\kappa$  KL 散度阈值)的影响。
- 3. **大规模验证实验**:最后,将最有潜力的设置应用到更大规模的训练中(1T tokens),以验证该方法在接近真实世界应用场景下的有效性和可扩展性。

## 4.1 合成数据集

### 原文翻译

To confirm that the Free Transformer indeed utilizes Z to condition its generative process, we designed a synthetic dataset and trained a small Free Transformer with different free-bits thresholds. Doing so allows to observe what aspects of the modeling are packed by the encoder in Z.

为了确认自由Transformer确实利用Z来调节其生成过程,我们设计了一个合成数据集,并用不同的自由比特阈值训练了一个小型的自由Transformer。这样做可以让我们观察到编码器将建模的哪些方面打包进了Z中。

### 深度解读

这个实验的目的非常明确:做一个"思想实验"的可视化验证。研究者想知道,当模型被赋予使用潜变量Z的能力时,它会优先用Z来记住哪些信息?通过控制"信息预算"*K*的大小,他们可以像调节阀门一样,观察在信息通道大小不同时,模型会做出什么样的权衡和取舍。

### 原文翻译

Each sequence in our synthetic training set is generated as follows:

- start with a string of 64 underscores "\_",
- pick an upper case letter and a position in the sequence at random, and replace the underscores there with a "target" made of the selected letter repeated 8 times,
- ullet replace any character with an exclamation mark with probability 1/16
- concatenate a prompt made of the target's letter followed by a ">".

我们合成训练集中的每个序列按如下方式生成:

• 以一个由64个下划线""组成的字符串开始,

- 随机选择一个大写字母和一个序列中的位置,然后用由该字母重复8次组成的"目标"替换那里的下划线,
- 以 1/16 的概率将任何字符替换为感叹号,
- 连接一个由目标字母和一个">"组成的提示符。

### 图3描述

这张图展示了合成数据的几个例子。我们可以清晰地看到数据的构成要素:

- 背景: 大量的下划线 。
- 目标: 在一个随机位置出现的、由同一个字母(如K, C, X)重复8次构成的字符串。
- 噪声: 随机散布的感叹号!。
- **提示**:在序列开头,用 字母〉的形式明确告知模型本次生成的目标字母是什么。

这个数据集设计得非常巧妙,因为它包含了不同层次的、需要模型去学习和重构的信息:

- 1. 最高层、最关键的信息:目标字符串出现的位置。
- 2. 次要信息: 随机噪声(感叹号)的位置和数量。
- 3. 已知信息:目标字母本身(已在提示中给出)。

模型的任务就是根据提示(如 K> ),补完整个64个字符的序列。一个好的模型需要准确地重构出目标的位置和噪声。

#### 原文翻译

We trained a Free Transformer on this data for four different values of the free bits threshold  $\kappa$ , and generated with the same random prompt three groups of sequences with each model, as pictured in Figure 4. For each model, in the blue group, the noise Z is sampled independently for each sequence, whereas we sampled one Z only for each of the green groups, used to generate all its sequences.

我们用这个数据针对四种不同的自由比特阈值  $\kappa$ 训练了自由Transformer,并用每个模型以相同的随机提示生成了三组序列,如图4所示。对于每个模型,在蓝色组中,噪声Z是为每个序列独立采样的;而在绿色组中,我们只为每个组采样一个Z,并用它来生成该组中的所有序列。

### 深度解读

这里的实验设置是理解图4结果的关键。

- 蓝色框:每次生成都用一个新的、随机的Z。这模拟了模型的正常使用方式。
- 绿色框:整个框内的所有生成,都使用同一个固定的Z。这是一个诊断性的实验。如果绿色框内的所有序列在某个方面表现出高度的一致性(比如目标都在同一个位置),那就强有力地证明了,这个一致的方面,正是被编码在那个固定的Z里面的信息。

#### 图4描述与分析

这张图是本次实验的核心结果,它完美地展示了模型的"思维"过程。让我们逐个分析:

- 左上角:  $\kappa = log(2)/64$  (1/64 bit)
  - 现象:信息预算极低。无论是蓝色框还是绿色框,模型生成的序列中,目标(TTTTTTTT)
    出现的位置都是随机的。绿色框内序列之间没有任何共性。
  - **结论**: 当信息预算几乎为零时,模型完全没有使用潜变量Z。它退化成了一个标准的自回归模型,每次生成都像是在"碰运气",无法规划目标的位置。
- 右上角:  $\kappa = log(2)/8$  (1/8 bit)
  - 现象:信息预算有所增加。蓝色框中,目标位置依然是随机的。但关键看绿色框:在第一个绿色框中,所有的目标(JJJJJJJJ)都出现在了序列的开头;在第二个绿色框中,所有的目标(FFFFFFF)都出现在了序列的中间。然而,每个序列中的噪声(感叹号)位置仍然是随机的。
  - **结论**: 这清晰地表明,当有了一点点信息预算后,模型学到了最重要的事情: **用潜变量Z来 编码目标的位置**。因为位置信息对于重构整个序列是最关键的。它还没有足够的"预算" 去记忆噪声的位置。
- 左下角: κ = log(2) (1 bit)
  - 现象:信息预算充足。现在,不仅绿色框内目标的位置是固定的,连噪声(感叹号)的位置也变得完全一致了。
  - **结论**:有了更多的信息预算,模型不仅用Z编码了目标的位置,还把随机噪声的位置和模式也一并编码了进去。此时,Z几乎包含了生成一个确定序列所需的所有"可变"信息。
- 右下角: κ = 8log(2) (8 bits)
  - 现象:信息预算过高。模型生成的内容开始出现错误和混乱。
  - **结论**: 这对应了我们之前讨论的"作弊"问题。当信息预算过于宽松时,编码器开始试图 将整个序列的原始信息(包括提示符本身)都硬塞进Z里,导致解码器变得"懒惰",不再 自己进行有效的生成。这破坏了VAE的训练平衡,导致性能下降。

### 综合洞察 这个合成实验有力地证明了:

- 1. Free Transformer确实能够利用潜变量Z来编码和控制生成过程的全局结构。
- 2. 模型在利用Z时表现出一种"智能"的优先级排序:它会优先编码对任务最重要的信息(位置),其次才是次要信息(噪声)。这体现了"信息瓶颈"原理的有效性。
- 3. "自由比特"阈值 K是一个关键的超参数,它需要被仔细调节。太小则Z无效,太大则模型 "作弊",两者都会损害性能。

### 4.2 基线架构

### 原文翻译

For assessing performance on standard benchmarks we used decoder-only Transformers implemented in the same Meta FAIR Transformer codebase as the one used by Copet et al. (2025) for the Computational World Model. Those are well optimized models using the

SwiGLU non-linearity (Shazeer, 2020), pre-normalization with RMSNorm (Zhang et al., 2019), Rotary Positional Embedding (ROPE, Su et al. 2021), and Group Query Attention (GQA, Ainslie et al. 2023). The vocabulary size is  $2^{17} \approx 130 k$ .

为了在标准基准上评估性能,我们使用了纯解码器Transformer作为基线。这些模型在Meta FAIR Transformer代码库中实现,与Copet等人(2025)用于计算世界模型的代码库相同。这些是经过良好优化的模型,使用了SwiGLU非线性激活函数(Shazeer, 2020)、使用RMSNorm的预归一化(Zhang et al., 2019)、旋转位置嵌入(ROPE, Su et al. 2021)和分组查询注意力(GQA, Ainslie et al. 2023)。词汇表大小为  $2^{17} \approx 130 \textit{k}$ 。

## 深度解读

这一段定义了实验的"控制组"。为了证明新架构的优越性,必须将它与当前最先进的标准架构进行公平比较。作者在此列出了一系列技术术语(SwiGLU, RMSNorm, RoPE, GQA),你不需要理解它们的具体细节,只需要知道这些都是当前构建高性能大语言模型的"标配"组件。作者强调他们使用的基线模型是"经过良好优化的",这很重要,因为它表明任何性能提升都不是因为基线模型太弱,而是因为新架构本身确实更强。

## 原文翻译

We used two sizes of models:

- A 1.5B model, with 28 layers, weight tying between the embeddings and the logit readout, model dimension 1536, 12 query heads, and 2 key-value heads. It is trained with 47B tokens, which requires 32 H100s for  $\approx 12$  hours.
- A 8B model with the structure of a Llama-3, which is 32 layers, model dimension 4096, 32 query heads, and 8 key-value heads. It is trained with 200B tokens which requires 256 H100s for ≈ 24 hours, or with 1T tokens, which takes 5 days.

## 我们使用了两种尺寸的模型:

- 一个15亿(1.5B)参数的模型,有28层,嵌入层和logit读出层权重绑定,模型维度1536,12 个查询头,2个键值头。它用470亿(47B)词元进行训练,需要32个H100 GPU训练约12小 时。
- 一个80亿(8B)参数的模型,其结构与Llama-3类似,有32层,模型维度4096,32个查询头,8个键值头。它用2000亿(200B)词元进行训练,需要256个H100 GPU训练约24小时;或者用1万亿(1T)词元训练,耗时5天。

## 深度解读

这里详细列出了实验所用模型的具体参数和训练成本。这不仅展示了研究的规模和投入,也为其他研究者复现实验提供了必要的信息。1.5B和8B是当前开源社区中非常流行的模型尺寸,覆盖了中等到较大的模型范围。训练成本的披露(例如256个H100 GPU训练24小时)也直观地体现了现代大模型研究的高昂门槛。

#### 原文翻译

We compare those baselines to the equivalent Free Transformers, which require one additional layer for the encoder during training and KV cache pre-filling, resulting in a compute and memory overhead of  $1/28 \approx 3.6\%$  for the 1.5B and  $1/32 \approx 3.1\%$  for the 8B.

我们将这些基线与等效的自由Transformer进行比较,后者在训练和KV缓存预填充时需要为编码器增加一个额外的层,导致1.5B模型的计算和内存开销增加  $1/28 \approx 3.6\%$ ,8B模型增加  $1/32 \approx 3.1\%$ 。

## 深度解读

这是对Free Transformer效率的量化说明,也是本文的一个核心卖点。增加的额外成本非常小:对于一个32层的模型,只增加了一层专用于编码器的模块,开销大约是3%。这意味着,可以用极小的代价,换取潜在的巨大性能提升。这使得该技术在商业上和学术上都极具吸引力。

## 4.3 设置和超参数

## 原文翻译

We kept our findings as clear as possible by avoiding other sources of performance improvement:

- We stuck to the baseline architecture, optimizer, and learning rate schedule that were used to train the baselines in FAIR's framework, and did not optimize any hyperparameter for our setup.
- We avoided any recipes for the VAE components, such as removing sampling in inference.
  We followed the formal expressions rigorously.
- We fixed H to 16 so that the dimension of  $Z_t$  was comparable to the vocabulary size of  $2^{17}$ .

我们通过避免其他性能提升来源,尽可能保持研究结果的清晰性:

- 我们坚持使用FAIR框架中用于训练基线的架构、优化器和学习率策略,并且没有为我们的设置 优化任何超参数。
- 我们避免了任何针对VAE组件的"秘方",例如在推理中移除采样。我们严格遵循了正式的数学表达式。
- 我们将H固定为16,使得  $Z_t$  的维度与  $2^{17}$  的词汇表大小具有可比性。

### 深度解读

这一段体现了严谨的科学精神。作者为了确保实验的公平性和结果的纯粹性,做出了以下承诺:

- 1. **不针对性调优**: 他们没有为Free Transformer "开小灶",而是沿用了为基线模型精心调校好的 所有训练参数。这意味着,如果Free Transformer表现更好,那完全是其架构本身的功劳,而 不是因为它享受了更好的训练待遇。这使得结果更加令人信服。
- 2. **忠于理论**: 他们没有使用一些可能会提升性能但缺乏理论依据的"黑科技"或"炼丹术",而是严格按照VAE的数学原理来实现模型。

3. **参数选择的合理解释**: 他们解释了为什么选择H=16,即让潜变量空间的维度( $2^{16}$ )与词汇表的大小( $2^{17}$ )在同一个数量级上,这是一个合理的设计选择。

### 原文翻译

We stress that the optimization hyperparameters were highly tuned for the baselines, and it is probable that a combination of an encoder and a decoder has specific requirements that would greatly benefit from an adapted training procedure.

我们强调,优化超参数是为基线模型高度调优的,而一个编码器和解码器的组合很可能有其特定的 需求,如果采用与之相适应的训练程序,可能会获得更大的收益。

### 深度解读

这是作者非常坦诚和自信的表态。他们等于在说:"我们是在一个'不公平'的比赛中获胜的。所有的训练设置都是为对手(基线模型)量身定做的。即便如此,我们的新模型还是表现得更好。如果未来有人专门为我们的Free Transformer设计一套最优的训练方案,它的性能潜力将会更大。" 这不仅增强了当前结果的说服力,也为未来的研究指明了方向。

## 4.4 探索性结果

## 原文翻译

We ran a series of experiments to assess the general behavior of the Free Transformer, and to calibrate the  $\kappa$ threshold. For any value of  $\kappa$ , the cross-entropy goes down regularly during training, with no more instability and spikes than what happens with the baselines. The KL divergence rapidly goes under  $\kappa$  and stays there.

我们进行了一系列实验来评估自由Transformer的总体行为,并校准  $\kappa$  阈值。对于任何  $\kappa$  值,交叉 熵在训练过程中都会稳定下降,其不稳定性和尖峰不比基线模型多。 $\kappa$  KL 散度会迅速下降到  $\kappa$  以下并保持在该水平。

### 深度解读

这是对训练过程稳定性的描述。结果表明,引入VAE结构并没有让训练过程变得更困难或不稳定,KL 散度也如预期一样被"自由比特"机制有效控制。这是一个积极的信号,说明新模型是"驯良" 的,易于训练。

### 原文翻译

When we compare the cross-entropies for various  $\kappa$ , they go down when  $\kappa$  increases as expected, but the values remain extremely close, with a difference of the order of 0.01 for a cross-entropy of  $\approx 2$  for the 1.5B and  $\approx 1.8$  for the 8B.

当我们比较不同 K值下的交叉熵时,正如预期的那样,随着 K的增加,交叉熵会下降,但这些值仍然非常接近,对于1.5B模型约2的交叉熵和8B模型约1.8的交叉熵,差异仅在0.01的量级。

### 深度解读

交叉熵是衡量模型预测下一个词准确性的指标,值越低越好。这里观察到一个有趣的现象:增加 "信息预算" *K*(允许Z携带更多信息)确实能让模型更好地拟合训练数据(交叉熵降低),但这种提升非常微小。这说明,仅仅在拟合训练数据这个层面,潜变量带来的帮助有限。真正的优势需要到更复杂的下游任务中去检验。

## 原文翻译

For both sizes of models, setting  $\kappa = 4 \log 2$ , corresponding to 4 bits of information per token, resulted in a collapse of the cross-entropy, indicating that the encoder found a way to channel fully the tokens to predict, and resulting in a collapse of performance on the downstream tasks. It is noteworthy that the baseline 8B model reaches during training a cross-entropy of  $1.8 = 2.59 \log(2)$ , hence may explain why allowing 2 bits does not collapse, while allowing 4 bits does.

对于两种尺寸的模型,当设置  $K = 4 \log 2$ (对应每个词元4比特信息)时,都导致了交叉熵的崩溃,这表明编码器找到了一种方法来完全传递要预测的词元,从而导致了下游任务性能的崩溃。值得注意的是,基线8B模型在训练中达到的交叉熵为  $1.8 = 2.59 \log(2)$ ,这或许可以解释为什么允许 2比特信息不会崩溃,而允许4比特信息则会。

## 深度解读

这验证了合成实验中的观察:过高的"信息预算"会导致模型"作弊"。当 K设置为4比特时,编码器有足够的带宽将下一个词的答案直接"泄露"给解码器,导致模型在训练集上表现完美(交叉熵崩溃),但在实际任务中一败涂地,因为它没有学到任何真正的生成能力。作者还敏锐地观察到,一个模型能从潜变量中有效利用的信息量,似乎与其自身的预测能力(由其最低交叉熵体现)有关。8B模型的预测不确定性大约是2.59比特,所以给它超过这个值的信息预算就可能导致"信息过载"和作弊。这是一个深刻的观察,揭示了模型容量、预测能力和潜变量信息量之间的内在联系。

## 原文翻译

The performance on downstream tasks are given in Table 1 for the 1.5B models, and Table 2 for the 8B models, both for four different values of  $\kappa$  corresponding to 1/2 to 2 bits of information per token. Graphs of performance during training are given in Appendix C in Figures 5 and 6. We observe a substantial increase of performance on HumanEval+, MBPP, and GSM8K which are arguably the benchmarks requiring some form of reasoning, and there also is a clear improvement for the 8B model with 1/2 bit of KL divergence on MMLU and CSQA, which are multi-choice questions.

1.5B模型在下游任务上的性能如表1所示,8B模型的性能如表2所示,两者都针对四种不同的 *K* 值(对应每个词元1/2到2比特的信息)。训练过程中的性能图表见附录C的图5和图6。我们观察到,在 HumanEval+、MBPP和GSM8K这些可以说需要某种形式推理的基准上,性能有显著提升;同时,对于8B模型,在使用1/2比特KL散度的设置下,在MMLU和CSQA这两个多项选择题基准上也有明显 改进。

### 深度解读

这是实验结果的核心摘要,也是全文最有力的证据。作者发现,Free Transformer的优势并非平均分布在所有任务上,而是集中在**需要推理能力**的任务上,具体包括:

- 代码生成 (HumanEval+, MBPP): 编写代码需要严密的逻辑推理和对程序结构的整体把握。
- 数学应用题 (GSM8K):解决数学问题需要遵循清晰的、多步骤的推理链条。

此外,在一些需要知识和常识的**问答任务**(MMLU, CSQA)上,新模型也表现出优势。 这个发现与我们之前的理论推导完美契合。潜变量Z扮演的"隐藏计划"或"高层大纲"角色,对于那些需要结构化、多步骤思考的任务来说,价值最大。对于简单的、记忆性的任务,这种"预先规划"可能就不那么重要了。这表明,Free Transformer不仅仅是提升了模型的性能,更是以一种有针对性的方式,增强了模型的"推理"能力。

表1: 1.5B模型(47B词元)性能

类 别	基准测试	基线	Free Transformer (1/4 bit)	Free Transformer (1/2 bit)	Free Transformer (1 bit)	Free Transforme (2 bits)
生成式代码数学	human_eval_plus (pass@1)	0.055	0.079 (+44.44%)	0.085 (+55.56%)	0.085 (+55.56%)	0.079 (+44.44%)
	mbpp (pass@1)	0.112	0.144 (+28.57%)	0.148 (+32.14%)	0.152 (+35.71%)	0.122 (+8.93%)
	gsm8k (em)	0.025	0.028 (+12.12%)	0.027 (+6.06%)	0.033 (+30.30%)	0.027 (+6.06%)
多选通用知识常识	mmlu (macro_avg/acc_char)	0.252	0.265 (+5.31%)	0.261 (+3.76%)	0.254 (+1.07%)	0.257 (+2.19%)
	csqa (acc_char)	0.199	0.199 (+0.00%)	0.187 (-6.17%)	0.197 (-0.82%)	0.175 (-11.93%)
	hellaswag (acc_char)	0.593	0.591 (-0.40%)	0.594 (+0.15%)	0.592 (-0.27%)	0.595 (+0.32%)
	winogrande (acc_char)	0.603	0.604 (+0.13%)	0.598 (-0.79%)	0.600 (-0.52%)	0.597 (-1.05%)
	obqa (acc_completion)	0.446	0.450 (+0.90%)	0.468 (+4.93%)	0.460 (+3.14%)	0.490 (+9.87%)

类 别	基准测试	基线	Free Transformer (1/4 bit)	Free Transformer (1/2 bit)	Free Transformer (1 bit)	Free Transforme (2 bits)
	arc_challenge (acc_completion)	0.400	0.392 (-1.93%)	0.386 (-3.43%)	0.405 (+1.29%)	0.385 (-3.65%)
	arc_easy (acc_completion)	0.596	0.602 (+0.92%)	0.592 (-0.64%)	0.603 (+1.06%)	0.592 (-0.71%)
	piqa (acc_char)	0.734	0.736 (+0.22%)	0.738 (+0.52%)	0.734 (+0.07%)	0.733 (-0.15%)
多选文本理解	race.high (acc_char)	0.390	0.382 (-2.20%)	0.390 (+0.00%)	0.387 (-0.81%)	0.386 (-1.03%)
	race.middle (acc_char)	0.532	0.511 (-3.93%)	0.519 (-2.49%)	0.522 (-1.83%)	0.514 (-3.40%)
	boolq (acc_completion)	0.583	0.632 (+8.39%)	0.614 (+5.35%)	0.648 (+11.12%)	0.620 (+6.29%)
文 化	nq (em)	0.081	0.069 (-15.36%)	0.073 (-9.56%)	0.075 (-7.17%)	0.071 (-11.95%)
	tqa (em)	0.205	0.191 (-6.93%)	0.190 (-7.58%)	0.200 (-2.84%)	0.197 (-4.13%)

## 表1解读 这张表格的数据非常震撼。在1.5B这个中等规模的模型上:

- **推理能力爆增**:在代码和数学任务上,Free Transformer取得了惊人的提升。HumanEval+提升超过55%,MBPP提升超过35%,GSM8K提升超过30%。这些都是两位数的、巨大的性能飞跃,在成熟的基准测试上是极为罕见的。
- 知识问答有提升:在MMLU和OBQA等需要知识和推理结合的任务上,也看到了不错的正向收益。
- **部分任务无变化或下降**:在一些常识推理(Hellaswag, Winogrande)和文本理解(RACE)任务上,性能基本持平。而在一些依赖精确事实抽取的任务(NQ, TQA)上,性能甚至有所下降。
- **最佳\$\kappa\$值**:不同的任务似乎偏好不同的"信息预算"。但总体来看,1/2比特到1比特似 平是一个比较理想的区间。

综合洞察 这些结果清晰地描绘了Free Transformer的"能力画像":它是一个专精于结构化推理的模型。对于需要逻辑、规划和多步骤思考的任务,潜变量Z提供的"全局大纲"能力起到了决定性的作用。而对于那些更侧重于从文本中查找和复述一个孤立事实的任务(比如NQ的"…的首都是哪里"),这种全局规划能力不仅帮助不大,甚至可能因为引入了额外的随机性而产生干扰。

## 表2: 8B模型 (200B词元) 性能

类 别	基准测试	基线	Free Transformer (1/4 bit)	Free Transformer (1/2 bit)	Free Transformer (1 bit)	Free Transforme (2 bits)
生成式代码数学	human_eval_plus (pass@1)	0.159	0.171 (+7.69%)	0.189 (+19.23%)	0.165 (+3.85%)	0.177 (+11.54%)
	mbpp (pass@1)	0.278	0.330 (+18.71%)	0.306 (+10.07%)	0.298 (+7.19%)	0.318 (+14.39%)
	gsm8k (em)	0.095	0.086 (-8.77%)	0.104 (+9.65%)	0.114 (+20.18%)	0.096 (-10.53%)
多选通用知识常识	mmlu (macro_avg/acc_char)	0.359	0.337 (-6.13%)	0.398 (+10.97%)	0.365 (+1.81%)	0.345 (-4.00%)
	csqa (acc_char)	0.356	0.292 (-17.93%)	0.450 (+26.21%)	0.346 (-2.99%)	0.324 (-8.97%)
	hellaswag (acc_char)	0.735	0.737 (+0.26%)	0.737 (+0.26%)	0.732 (-0.45%)	0.738 (+0.39%)
	winogrande (acc_char)	0.680	0.667 (-1.86%)	0.664 (-2.32%)	0.664 (-2.32%)	0.667 (-1.86%)
	obqa (acc_completion)	0.522	0.508 (-2.68%)	0.484 (-7.28%)	0.530 (+1.53%)	0.554 (+6.13%)
	arc_challenge (acc_completion)	0.465	0.483 (+3.87%)	0.468 (+0.55%)	0.452 (-2.95%)	0.485 (+4.24%)
	arc_easy (acc_completion)	0.677	0.676 (-0.25%)	0.665 (-1.81%)	0.668 (-1.44%)	0.679 (+0.31%)

类 别	基准测试	基线	Free Transformer (1/4 bit)	Free Transformer (1/2 bit)	Free Transformer (1 bit)	Free Transforme (2 bits)
	piqa (acc_char)	0.774	0.780 (+0.77%)	0.782 (+1.05%)	0.785 (+1.41%)	0.793 (+2.46%)
多选文本理解	race.high (acc_char)	0.433	0.447 (+3.30%)	0.443 (+2.25%)	0.444 (+2.58%)	0.435 (+0.53%)
	race.middle (acc_char)	0.594	0.592 (-0.35%)	0.591 (-0.47%)	0.587 (-1.17%)	0.584 (-1.64%)
	boolq (acc_completion)	0.705	0.632 (-10.37%)	0.632 (-10.33%)	0.687 (-2.47%)	0.671 (-4.82%)
文 化	nq (em)	0.181	0.183 (+1.38%)	0.167 (-7.67%)	0.173 (-4.14%)	0.168 (-6.90%)
	tqa (em)	0.438	0.440 (+0.28%)	0.443 (+0.80%)	0.434 (-1.19%)	0.446 (+1.45%)

## 表2解读 在更大的8B模型上,我们看到了类似的趋势,但有一些新的变化:

- **推理能力持续领先**:在代码和数学任务上,Free Transformer继续保持着两位数的巨大优势 (例如HumanEval+提升19%,GSM8K提升20%)。
- 知识问答优势扩大:在MMLU和CSQA上,提升幅度变得更加惊人,分别达到了11%和26%。这表明随着模型本身知识量的增加,潜变量提供的"推理框架"能更有效地组织和运用这些知识。
- $\kappa = 1/2$  **bit 表现突出**:在许多任务中, $\kappa = 1/2$  bit 这个设置取得了最好的或接近最好的结果。 这似乎是8B模型的一个"甜点"区域。
- 任务表现波动:与1.5B模型相比,不同任务上的表现波动更大。例如,在BoolQ上性能显著下降,但在RACE.high上则有提升。这可能说明,新架构与不同任务的契合度,以及最佳超参数的选择,可能与模型规模有关,需要更精细的调优。

## 4.5 1万亿词元训练结果

### 原文翻译

To measure improvement in a more realistic setting, closer to models actually used in real applications, we trained 8B models on 1T tokens, which improves drastically the performance of both the baseline and the Free Transformer. Given the results with 200B tokens, we chose the value  $\kappa = log(2)/2$  corresponding to half a bit of information per token at most.

为了在更现实的、更接近实际应用模型的环境中衡量改进,我们在1万亿(1T)词元上训练了8B模型,这极大地提升了基线和自由Transformer的性能。根据200B词元训练的结果,我们选择了  $\kappa = \log(2)/2$  的值,对应每个词元最多半比特的信息。

## 深度解读

这是最终的、也是最具说服力的实验。训练数据量从2000亿增加到1万亿,这是一个巨大的飞跃,能让模型学到远比之前更丰富和深入的知识与能力。作者根据之前的经验,明智地选择了表现最稳健的  $\kappa = 1/2$  bit 作为本次实验的设置。这次实验的目的是检验:在模型本身已经非常强大的情况下,Free Transformer带来的"结构化推理"优势是否依然存在,甚至被放大?

## 原文翻译

The performance on downstream tasks are given in Table 3 and the corresponding graphs during training in Figure 7 of Appendix C. We provide in the table the performance measured at the end of the training as for the other configurations, but in addition we also give the average over the last third of the training. We can observe on the graphs that the rate of improvement tend to be constant on this interval, which justifies averaging to mitigate the performance fluctuations.

下游任务的性能在表3中给出,相应的训练过程图表在附录C的图7中。我们在表中提供了训练结束时测量的性能,与其他配置一样;但此外,我们还给出了训练最后三分之一阶段的平均性能。我们可以从图表中观察到,在这个区间内,性能提升的速率趋于稳定,这为通过平均来减轻性能波动提供了理由。

## 深度解读

作者在这里增加了一个更稳健的评估指标:最后三分之一训练周期的平均性能。因为模型在训练末期的性能可能会有随机波动,只看最后一个点可能会有偶然性。取一段时间的平均值可以更公允地 反映模型的真实水平。

### 原文翻译

The key result is the boost of performance on HumanEval+, MBPP, GSM8K, MMLU and CSQA, confirming what we observed in the smaller settings, and a greater stability on other tasks.

关键结果是在HumanEval+、MBPP、GSM8K、MMLU和CSQA上的性能提升,这证实了我们在较小规模设置中观察到的现象,并且在其他任务上表现出更大的稳定性。

## 深度解读

一言以蔽之:之前的优势得到了保持和确认。即使在1T tokens的强力训练下,Free Transformer依然在代码、数学和知识问答这些核心推理任务上,展现出比同样强大的基线模型更胜一筹的能力。同时,在其他任务上的表现也变得更加稳定,之前看到的一些性能下降现象有所缓解。这表明,Free Transformer带来的改进是根本性的,并且随着模型能力的增强而"水涨船高"。

表3:8B模型(1T词元)性能

类别	基准测试	基线 (最 终值)	FT 1/2 bit (最 终值)	基线 (后1/3 均值)	FT 1/2 bit (后 1/3均值)
生成式代 码/数学	human_eval_plus (pass@1)	0.268	0.299 (+11.36%)	0.245	0.256 (+4.22%)
	mbpp (pass@1)	0.428	0.440 (+2.80%)	0.396	0.421 (+6.08%)
	gsm8k (em)	0.321	0.331 (+2.83%)	0.280	0.296 (+5.84%)
多选-通用 知识/常识	mmlu (macro_avg/acc_char)	0.592	0.623 (+5.20%)	0.567	0.596 (+5.16%)
	csqa (acc_char)	0.707	0.748 (+5.79%)	0.689	0.733 (+6.28%)
	hellaswag (acc_char)	0.799	0.799 (-0.01%)	0.787	0.788 (+0.18%
	winogrande (acc_char)	0.739	0.735 (-0.53%)	0.725	0.727 (+0.27%
	obqa (acc_completion)	0.564	0.562 (-0.35%)	0.556	0.551 (-0.86%)
	arc_challenge (acc_completion)	0.542	0.535 (-1.42%)	0.524	0.522 (-0.40%)
	arc_easy (acc_completion)	0.721	0.711 (-1.41%)	0.706	0.711 (+0.68%
	piqa (acc_char)	0.805	0.812 (+0.88%)	0.802	0.807 (+0.61%)
多选-文本 理解	race.high (acc_char)	0.473	0.463 (-2.06%)	0.467	0.460 (-1.55%)
	race.middle (acc_char)	0.632	0.634 (+0.33%)	0.623	0.624 (+0.16%
	boolq (acc_completion)	0.713	0.725 (+1.63%)	0.755	0.754 (-0.10%)
文化	nq (em)	0.248	0.247 (-0.22%)	0.229	0.227 (-0.76%)

类别	基准测试	基线 (最 终值)	FT 1/2 bit (最 终值)	基线 (后1/3 均值)	FT 1/2 bit (后 1/3均值)
	tqa (em)	0.583	0.577 (-1.00%)	0.549	0.544 (-0.90%)

### 表3解读 这张表格的数据证实了作者的结论:

- **核心优势稳固**:在HumanEval+, MBPP, GSM8K, MMLU, CSQA这五大核心推理基准上,Free Transformer(FT)无论是看最终值还是看平均值,都稳定地优于基线模型。提升幅度虽然不像在小模型上那么夸张(因为基线本身已经很强),但依然是显著的、有意义的进步(普遍在2%到11%之间)。
- **稳定性提升**:在Hellaswag, Winogrande等之前表现不佳或持平的任务上,性能差距已经缩小到几乎可以忽略不计。这表明更大规模的训练帮助Free Transformer更好地学习了如何运用潜变量,减少了它在不擅长任务上的负面影响。
- **均值指标的价值**: 比较"最终值"和"后1/3均值"两列,我们可以看到,基于平均值的相对提升(如MBPP +6.08%, GSM8K +5.84%)甚至比基于最终值的提升更明显。这印证了作者的说法,即平均值能更好地反映Free Transformer在训练后期稳定超越基线的趋势。

**实验部分总结** 从一个精心设计的合成实验,到两个尺寸、两种训练规模的全面基准测试,实验部分以无可辩驳的数据证明了Free Transformer的核心价值:通过引入一个低成本的、由VAE框架训练的潜变量,模型能够学习到一种"先规划后执行"的生成模式,从而显著提升其在需要结构化推理的任务上的表现。

# 5. 相关工作

### 原文翻译

There have been several attempts at combining a VAE and a decoder Transformer, generally with a focus on improving topic models and providing ways to guide the generation. The OPTIMUS model (Li et al., 2020) combines a pre-trained BERT as text embedding/encoder, with a GPT-2 playing the role of decoder, which are fine-tuned with a VAE-like loss. The latent embedding Z is computed thanks to a CLS token, that is by adding a token to the input and a read-out to extract its embedding in the output. To modulate the GPT-2 generation with it, it is either (1) concatenated as an additional token in every layer, or (2) added to the input token embeddings. Collapse of the KL divergence is prevented during training with the free bits method (Kingma et al., 2016). This approach allows for better guided text generation with GPT-2 and better generalization on low-data languages with BERT.

已经有数次将VAE与解码器Transformer结合的尝试,通常侧重于改进主题模型和提供引导生成的方法。OPTIMUS模型(Li et al., 2020)将一个预训练的BERT作为文本嵌入/编码器,与一个扮演解码器角色的GPT-2相结合,并通过类似VAE的损失进行微调。潜变量嵌入Z是通过一个CLS词元计算得到的,即在输入中增加一个词元,并通过一个读出层来提取其在输出中的嵌入。为了用它来调节GPT-2的生成,它要么(1)作为额外的词元在每一层进行拼接,要么(2)加到输入词元嵌入上。训练中通过自由比特方法(Kingma et al., 2016)来防止KL散度的坍塌。这种方法使得GPT-2能进行更好的引导式文本生成,并让BERT在低资源语言上具有更好的泛化能力。

### 深度解读

这一段回顾了前人的工作,将本研究置于更广阔的学术背景中。科学研究不是凭空产生的,而是在前人工作的基础上不断演进。作者在这里展示了他们对该领域的了解。 OPTIMUS模型是早期的一个重要尝试。但它的做法是"拼接": 把一个现成的编码器模型(BERT)和一个现成的解码器模型(GPT-2)硬凑在一起。这种方法的缺点是:

- 1. 结构笨重:需要维护两个独立的、巨大的模型。
- 2. **融合粗糙**:将潜变量Z注入解码器的方式比较简单,要么作为第一个词元,要么直接加在输入层,可能无法在解码器的深层结构中发挥充分作用。 尽管如此,OPTIMUS证明了这条路是可行的,并验证了"自由比特"等技巧的有效性。

## 原文翻译

Xie et al. (2021) extend OPTIMUS with a multi-objective loss, adding in particular the prediction of the story topic, using the output of another model as ground truth, to obtain a better embedding space.

Xie等人(2021)通过一个多目标损失扩展了OPTIMUS,特别是增加了对故事主题的预测,使用另一个模型的输出作为真实标签,以获得更好的嵌入空间。

## 深度解读

这项工作试图通过增加额外的、有监督的任务(如预测主题)来让潜变量空间Z的结构性更强。这是一种"半监督"的思路,但它依赖于外部模型来提供"主题"标签,增加了系统的复杂性。

## 原文翻译

The CVAE proposed by Fang et al. (2021) combines two pre-trained GPT-2, one used as the encoder without causal masking. The embedding Z is an average of the encoder's output, and the authors propose three ways to modulate the decoder with linear images of it: (1) add it to each input token embedding, (2) concatenate it to the Ks and Vs in every layer, (3) add it before the softmax. Experiments demonstrate that this method allows controlling the generation without hurting the quality of the result.

Fang等人(2021)提出的CVAE结合了两个预训练的GPT-2,其中一个在没有因果掩码的情况下用作编码器。嵌入Z是编码器输出的平均值,作者提出了三种用其线性变换来调节解码器的方法:(1)将其加到每个输入词元嵌入上,(2)将其拼接到每一层的K和V上,(3)在softmax之前添加它。实验表明,该方法可以在不损害结果质量的情况下控制生成。

## 原文翻译

AdaVAE (Tu et al., 2022) is similarly the combination of two pre-trained GPT-2, the first without causal masking playing the role of the encoder. The latent embedding Z is extracted from its output with a slightly modified attention operator. It is then injected into the decoder by either concatenating an image of it to the keys and values as in OPTIMUS, or before the softmax as in CVAE.

AdaVAE(Tu et al., 2022)类似地也是两个预训练GPT-2的组合,第一个在没有因果掩码的情况下扮演编码器的角色。潜变量嵌入Z是通过一个轻微修改的注意力算子从其输出中提取的。然后通过将其变换拼接到键和值上(如OPTIMUS),或在softmax之前(如CVAE)注入到解码器中。

### 深度解读

CVAE和AdaVAE都延续了使用两个独立GPT-2模型的思路,一个做编码器,一个做解码器。它们探索了更多将潜变量Z注入解码器的方式,比如在每一层的注意力机制中融入Z的信息。这些工作都为"可控生成"做出了贡献,但它们共同的局限性在于,都是基于"拼接"现有模型的思路,没有从根本上设计一个原生集成的、高效的VAE-Transformer架构。

与Free Transformer的对比 通过回顾这些工作,Free Transformer的创新之处就更加凸显了:

- **原生集成,而非拼接**:它不是把两个模型粘在一起,而是通过巧妙的层共享,将编码器"内生" 干解码器之中。
- 高效:由于层共享,其额外的计算和参数开销极小,这使得它可以被轻松地扩展到非常大的模型(如8B甚至更大),而之前的方法很难做到这一点。
- **无监督**:它完全依靠VAE的内在机制来学习有意义的潜变量,不需要像Xie等人的工作那样依赖 外部的监督信号。

# 6. 结论

## 原文翻译

The Free Transformer is a direct extension of a standard decoder Transformer, with the abstract structure of a conditional VAE. It is implemented with a single additional non-causal Transformer block and requires a few percent of computational and memory usage overhead. Its structure makes it able to learn latent random variables unsupervised, and to condition its generative process on them.

自由Transformer是标准解码器Transformer的一个直接扩展,具有条件VAE的抽象结构。它通过一个额外的非因果Transformer模块实现,只需要百分之几的计算和内存使用开销。其结构使其能够无监督地学习潜在随机变量,并以它们为条件来调节其生成过程。

#### 深度解读

这是对全文工作的凝练总结。作者重申了Free Transformer的几个关键特性:它是一个扩展,而非颠覆;它在理论上是一个cVAE;它在实现上高效;它在功能上实现了无监督的潜在空间学习和条件生成。

## 原文翻译

In some ways, this approach aims at achieving in latent space with an autoencoder what reasoning models do with chains-of-thought in token space and an RL procedure (DeepSeek-AI et al., 2025). A combination of the two is, of course, promising.

在某些方面,这种方法旨在通过自编码器在潜空间中实现的目标,与推理模型通过思维链(chainsof-thought)在词元空间中和通过强化学习(RL)程序实现的目标是相似的(DeepSeek-Al et al., 2025)。当然,将两者结合起来是很有前景的。

## 深度解读

这是一个极其深刻的洞察,它将Free Transformer与当前主流的另一种提升模型推理能力的技术——思维链(Chain-of-Thought, CoT)——进行了类比和联系。

- **思维链(CoT)**: 让模型在回答问题前,先把解题步骤一步步地写出来。这就像学生在做数学题时,被要求写出详细的解题过程。这种方式是"**显式的**"、"**在词元空间中进行的**"推理。模型是在"大声思考"。
- Free Transformer:模型在生成最终答案前,先在内部形成一个抽象的、高维的"计划"(潜变量Z)。这个计划是"**隐式的**"、"**在潜空间中进行的**"推理。模型是在"静默思考"。

作者认为,这两种方式殊途同归,都是为了给模型的生成过程引入结构和规划。CoT的优点是过程透明、可解释,但缺点是会增加输出的长度,且依赖于特定的提示技巧。Free Transformer的优点是高效、紧凑,但缺点是其"思考"过程是一个难以解释的"黑箱"。最后,作者提出了一个激动人心的展望:"将两者结合起来"。未来的模型或许可以同时拥有这两种能力:既能在潜空间中进行高效的、全局性的静默规划,又能将关键的推理步骤用自然语言"大声"地表达出来,以增强结果的可靠性和可解释性。这将是迈向更强大通用人工智能的重要一步。

## 原文翻译

The performance boost without tuning the optimization hyperparameters across multiple benchmarks and two sizes of models, is a strong signal that the overall approach actually improves the inductive bias of the vanilla Transformer.

在没有调整优化超参数的情况下,模型在多个基准和两种尺寸上都获得了性能提升,这是一个强烈的信号,表明该方法确实改善了原始Transformer的归纳偏置。

### 深度解读

这里引入了一个重要概念:"归纳偏置"(inductive bias)。你可以把它理解为一个学习模型的"世界观"或"先天假设"。

• 标准Transformer的归纳偏置:它假设文本的结构是纯粹的、局部的、序列性的。

• **Free Transformer的归纳偏置**:它假设文本背后存在一个隐藏的、全局的、层次化的结构。实验结果表明,后一种"世界观"更接近真实世界文本(尤其是需要推理的文本)的本质,因此模型能学得更好。在没有"应试"优化的情况下取得好成绩,说明这个学生(模型)是真的"天赋好"(归纳偏置好),而不是靠刷题技巧。

## 原文翻译

Many properties and design choices should be explored. The performance curves during training are often unstable, possibly due to the coupling of the optimization of the encoder and the decoder, and using different optimization methods could be fruitful. The random embedding itself could take many forms, and the one used in our implementation is arbitrary. Finally, the behavior in larger scales, both in parameter count and dataset size, remains to be investigated.

许多属性和设计选择有待探索。训练过程中的性能曲线通常不稳定,这可能是由于编码器和解码器 的优化耦合在一起,使用不同的优化方法可能会有成效。随机嵌入本身可以有多种形式,我们实现 中使用的是任意选择的一种。最后,其在更大规模(无论是参数数量还是数据集大小)下的行为仍 有待研究。

## 深度解读

这是对未来工作的展望,也体现了研究者的诚实。他们承认当前的工作还不是完美的:

- **训练不稳定性**: 同时优化编码器和解码器有时会导致训练过程出现波动。这可能需要新的、更 先进的优化算法来解决。
- **设计选择的探索**:潜变量Z的结构(比如为什么是65536维的独热向量)还有很大的探索空间,不同的结构可能适用于不同的任务。
- **规模扩展**:虽然已经在1T词元上进行了实验,但要真正了解其潜力,还需要在更大(如百亿、 千亿参数)的模型和更大(如十万亿、百万亿词元)的数据集上进行验证。

这篇论文为我们打开了一扇通往新型生成模型架构的大门。它以清晰的动机、优雅的理论、巧妙的 实现和详实的数据,证明了让模型在"静默中思考"的巨大潜力。

# 附录A 评估基准

#### 原文翻译与解读

- **HellaSwag**: Multiple choices. Common sense focusing on physically situated scenarios. (Zellers et al., 2019)
  - HellaSwag: 多项选择题。专注于物理场景下的常识推理。
  - 解读:测试模型对日常物理世界情景的理解。例如,给出"他打开了水龙头,"的开头,模型需要从几个选项中选出最合理的结尾,如"水流了出来"。
- **WinoGrande**: Large-scale adversarial Winograd-style pronoun resolution (fill-in-the-blank) designed to reduce annotation artifacts. (Sakaguchi et al., 2019)

- WinoGrande: 大规模对抗性的维诺格拉德模式代词消解(填空题),旨在减少标注偏差。
- **解读**:测试模型的指代消解和因果推理能力。例如,"奖杯放不进棕色的手提箱,因为它太大了。"模型需要判断"它"指的是奖杯还是手提箱。
- ARC (Al2 Reasoning Challenge): Grade-school science multiple choice. (Clark et al., 2018)
  - ARC (AI2推理挑战): 小学科学多项选择题。
  - **解读**:测试模型对基础科学知识的掌握和应用能力。
- **PIQA**: Physical commonsense multiple choice about everyday goals and affordances. (Bisk et al., 2019)
  - PIQA: 关于日常目标和物体功能的物理常识多项选择题。
  - 解读:测试模型对物体如何被使用的常识理解。例如,给出两个解决方案,模型需要选择哪一个更能达到目标。
- **OpenBookQA (OBQA)**: Open-book science QA: combines a provided set of core facts with commonsense/world knowledge to answer questions. (Mihaylov et al., 2018)
  - OpenBookQA (OBQA): 开卷科学问答: 结合一组给定的核心事实和常识/世界知识来回答问题。
  - 解**读**:测试模型结合给定信息和自身知识进行推理的能力,类似于开卷考试。
- **RACE**: Multiple-choice reading comprehension from Chinese middle-school English exams. (Lai et al., 2017)
  - RACE:来自中国初高中英语考试的多项选择阅读理解。
  - 解读:测试模型的长文本阅读和信息提取能力。
- MMLU: "Massive Multitask Language Understanding". Questions spanning STEM, humanities, social sciences, etc. (Hendrycks et al., 2021)
  - MMLU: "大规模多任务语言理解"。问题涵盖STEM、人文学科、社会科学等领域。
  - **解读**:这是一个综合性的知识水平测试,被认为是衡量模型综合智力的"高考"。
- **CommonsenseQA (CSQA)**: Multiple-choice QA requiring commonsense relational knowledge (leveraging ConceptNet relations). (Talmor et al., 2019)
  - CommonsenseQA (CSQA):需要常识性关系知识的多项选择问答(利用ConceptNet关系)。
  - **解读**:测试模型对概念之间常识性关系的理解。例如,"去哪里可以找到很多书?"答案是"图书馆"。
- **BoolQ**: Yes/no questions paired with passages to evaluate reading comprehension and entailment-like inference. (Clark et al., 2019)
  - BoolQ: 是/否问题,配有段落,用于评估阅读理解和类似蕴含的推理。
  - 解读:模型需要阅读一段文字,然后回答一个只能用"是"或"否"来回答的问题。

- **GSM8K**: Grade-school math word problems requiring multi-step arithmetic reasoning. (Cobbe et al., 2021)
  - **GSM8K**:需要多步算术推理的小学数学应用题。
  - **解读**:测试模型的数学逻辑推理和计算能力,是衡量模型推理能力的核心指标之一。
- **HumanEval+**: An augmented version of OpenAI's HumanEval (Chen et al., 2021) with many more unit tests per problem to reduce test fragility and overfitting in code generation evaluation. (Liu et al., 2023)
  - HumanEval+: OpenAl HumanEval的增强版,每个问题有更多的单元测试,以减少代码 生成评估中的测试脆弱性和过拟合。
  - 解读:测试模型根据自然语言描述生成正确代码的能力,是衡量代码能力的核心指标。
- MBPP: "Mostly Basic Programming Problems." Short Python programming tasks solvable by entry-level programmers; includes text spec and example tests. (Austin et al., 2021)
  - MBPP: "主要是基础编程问题"。入门级程序员可解决的简短Python编程任务;包括文本描述和示例测试。
  - 解读:另一个衡量基础代码生成能力的基准。
- **NQ**: "Natural Questions." Real user queries paired with Wikipedia pages. (Kwiatkowski et al., 2019)
  - NQ: "自然问题"。真实的用户查询与维基百科页面配对。
  - **解读**:测试模型从给定文档中精确查找答案的能力。
- **TQA**: "TriviaQA". A challenging reading comprehension dataset.
  - TQA: "知识问答"。一个具有挑战性的阅读理解数据集。
  - 解读: 类似于NQ,测试事实查找能力。

## 附录B 性能度量

### 原文翻译与解读

- For generated answers:
  - **pass@1** is the proportion of generated pieces of code that produce the expected behavior when executed.
    - pass@1: 生成的代码在执行时能够产生预期行为的比例。

我用这种等式的形式对这篇文章有如下的理解,你看看有什么要修正的 free\_transformer=潜变量Z的训练和应用=条件变分自编码器+全局性提问+自由比特超 参数控制=小模型的提升暴涨但不=大模型明显提升=与CoT相对应的推理模式构建

用户的理解基本准确,捕捉到了Free Transformer的核心机制和优势,但有一些细微之处需要修正和补充,以更精确地反映论文内容。 核心概念

Free Transformer确实是通过潜变量Z的训练和应用来扩展标准的decoder Transformer模型,其中Z是随机潜变量,通过无监督的变分过程学习。 这一过程等价于条件变分自编码器(cVAE)的框架,其中编码器(encoder)从序列中推断Z,解码器(decoder)则在生成时条件于Z,以实现更丰富的密度建模。 论文强调,这种条件化避免了纯自回归模型在隐式推断潜变量时的复杂性和错误累积。

关键组件条件变分自编码器:正确,这是基础架构。Z被注入Transformer的中层,通过一个非因果的encoder块来捕捉序列的全局属性(global properties),而非逐token映射。 这与用户的"全局性提问"相符,指的是encoder使用学习到的常量查询(learned constant input)来提取序列的全局结构,便于跨任务转移。自由比特超参数控制:精确描述了使用free bits方法来调控KL散度(Kullback-Leibler divergence),防止后验崩溃(posterior collapse)。 超参数阈值(如1-2 bits per token)控制Z中流动的信息量,避免encoder过度泄露序列细节给decoder。 论文中H=16(对应65,536维one-hot Z),通过binary mapper实现高效采样。

模型性能提升小模型提升暴涨:1.5B参数模型在47B tokens训练下,在推理任务(如HumanEval、MBPP、GSM8K)上提升显著,例如HumanEval从44.44%升至55.56%(+11.12%),GSM8K从6.06%至30.30%(+24.24%)。 这得益于Z提供明确的"决策点",类似于在潜空间中进行结构化生成。大模型提升:并非"不明显",8B参数模型也在200B或1T tokens训练下显示出明显改进,例如在1T tokens下HumanEval +11.36%、MMLU +5.20%、CSQA +5.79%。 不过,相对小模型的"暴涨",大模型的绝对提升更稳定但幅度稍小(例如GSM8K +2.83%),可能因基线已较强。 总体上,Free Transformer在多基准中提升推理能力,而在常识任务(如HellaSwag)上变化较小。与CoT的对应

用户的表述"与CoT相对应的推理模式构建"非常贴切。论文明确指出,Free Transformer在潜空间通过自编码器实现类似于Chain-of-Thought (CoT)在token空间的推理效果,后者依赖RL和显式链式提示。 Z允许模型在生成前做出全局"决策"(如正面/负面评论),类似于CoT的逐步推理,但更高效且无监督。 未来可结合两者,进一步提升大模型的泛化。

总体而言,你的等式框架简洁有力,仅需调整大模型提升的描述为"显著但相对稳定",以避免低估其效果。论文实验显示,这种设计在计算开销仅增加3-4%的情况下,改善了Transformer的归纳偏差(inductive bias)。The-Free-Transformer-Lun-Wen-Wan-Zheng-Fan-Yi-Yu-Shen-Du-Jie-Du.pdf?https://ppl-ai-file-upload.s3. amazonaws.com/web/direct-files/attachments/50767043/844c9813-c26f-4163-8154-4d6b80769cec/The-Free-Transformer-Lun-Wen-Wan-Zheng-Fan-Yi-Yu-Shen-Du-Jie-Du.pdf